

# Toward optimal model averaging in regression models with time series errors

Tzu-Chang F. Cheng, Ching-Kang Ing,  
Shu-Hui Yu

University of Illinois at Urbana-Champaign,  
Academia Sinica and National Taiwan University,  
National University of Kaohsiung

## Abstract

Consider a regression model with infinitely many parameters and time series errors. We are interested in choosing weights for averaging across generalized least squares (GLS) estimators obtained from a set of approximating models. However, GLS estimators, depending on the unknown inverse covariance matrix of the errors, are usually infeasible. We therefore construct feasible generalized least squares (FGLS) estimators using a consistent estimator of the unknown inverse matrix. Based on this inverse covariance matrix estimator and FGLS estimators, we develop a feasible autocovariance-corrected Mallows model averaging criterion to select weights, thereby providing an FGLS model averaging estimator of the true regression function. We show that the generalized squared error loss of our averaging estimator is asymptotically equivalent to the minimum one among those of GLS model averaging estimators with the weight vectors belonging to a continuous set, which includes the discrete weight set used in Hansen (2007) as its proper subset.

*JEL classification:* C22; C52

**KEY WORDS:** Asymptotic efficiency; Autocovariance-corrected Mallows model averaging; Banded Cholesky factorization; Feasible generalized least squares estimator; High-dimensional covariance matrix; Time series errors.

## 1 Introduction

This article is concerned with the implementation of model averaging methods in regression models with time series errors. We are interested in choosing weights for averaging across generalized least squares (GLS) estimators obtained from a set of approximating models for the true regression function. However, GLS estimators, depending on the unknown covariance matrix  $\Sigma_n^{-1}$  of the errors, are usually infeasible, where  $n$  is the sample size. We therefore construct feasible generalized least squares (FGLS) estimators using a consistent estimator of  $\Sigma_n^{-1}$ . Based on this inverse covariance matrix estimator and FGLS estimators, we develop a feasible autocovariance-corrected Mallows model averaging (FAMMA) criterion to select weights, thereby providing an FGLS model averaging estimator of the regression function. We show that the generalized squared error loss of our averaging estimator is asymptotically equivalent to the minimum one among those of GLS model averaging

estimators with the weight vectors belonging to a continuous set, which includes the discrete weight set used in Hansen (2007) as its proper subset.

Let  $M$  be the number of approximating models. If the weight set only contains standard unit vectors in  $R^M$ , then selection of weights for model averaging is equivalent to selection of models. Therefore, model selection can be viewed as a special case of model averaging. It is shown in Hansen (2007, p.1179) that when the weight set is rich enough, the optimal model averaging estimator usually outperforms the one obtained from the optimal single model, providing ample reason to conduct model averaging. Another vivid example demonstrating the advantage of model averaging over model selection is given by Yang (2007, Section 6.2.1, Figure 5). In the case of independent errors, asymptotic efficiency results for model selection have been reported extensively, even when the errors are heteroskedastic or regression functions are serially correlated. For the regression model with i.i.d. Gaussian errors, Shibata (1981) showed that Mallows'  $C_p$  (Mallows (1973)) and Akaike information criterion (AIC; Akaike (1974)) lead to asymptotically efficient estimators of the regression function. By making use of Whittle's (1960) moment bounds for quadratic forms in independent variables, Li (1987) established the asymptotic efficiency of Mallows'  $C_p$  under much weaker assumptions on homogeneous errors. Li's (1987) result was subsequently extended by Andrews (1991) to heteroscedastic errors. There are also asymptotic efficiency results established in situation where regression functions are serially correlated. Assuming that the data are generated from an infinite order autoregressive (AR( $\infty$ )) process driven by i.i.d. Gaussian noise, Shibata (1980) showed that AIC is asymptotically efficient for independent-realization prediction. This result was extended to non-Gaussian AR( $\infty$ ) processes by Lee and Karagrigoriou (2001). Ing and Wei (2005) showed that AIC is also asymptotically efficient for same-realization prediction. Ing (2007) further pointed out that the same property holds for a modification of Rissanen's accumulated prediction error (APE, Rissanen (1986)) criterion.

Asymptotic efficiency results for model averaging have also attracted much recent attention from econometricians and statisticians. Hansen (2007) proposed the Mallows model averaging (MMA) criterion, which selects weights for averaging across LS estimators. Under regression models with i.i.d. explanatory vectors and errors, he proved that the averaging estimator obtained from the MMA criterion asymptotically attains the minimum squared error loss among those of the LS model averaging estimators with the weight vectors contained in a discrete set  $\mathcal{H}_n(N)$  (see (2.8)), in which  $N$  is a positive integer and related to the moment restrictions of the errors. Using the same weight set, Hansen and Racine (2012) and Liu and Okui (2013), respectively, showed that the Jackknife model averaging (JMA) criterion and feasible  $HRC_p$  criterion yield asymptotically efficient LS model averaging estimators in regression models with independent explanatory vectors and heteroscedastic errors. Since  $\mathcal{H}_n(N)$  is quite restrictive when  $N$  is small, Wan, Zhang and Zou (2010) justified MMA's asymptotic efficiency over the continuous weight set

$$\mathcal{G}_n = \{\mathbf{w} = (w_1, \dots, w_M)' : w_m \in [0, 1], \sum_{m=1}^M w_m = 1\}, \quad (1.1)$$

which is much more flexible than  $\mathcal{H}_n(N)$ . Recently, Ando and Li (2014) showed that Hansen and Racine's (2012) result carries over to high-dimensional regression models and to a weight set more general than  $\mathcal{G}_n$ .

There are different types of theoretical examinations on model averaging. Besides the approach of targeting asymptotic efficiency, another very successful approach is minimax optimal model combination via oracle inequalities; see, for example, Yang (2001), Yuan and Yang (2005), Leung

and Barron (2006), and Wang et al. (2014).

However, all aforementioned papers, requiring the error terms to be independent, preclude the regression model with time series errors, which is one of the most useful models for analyzing dependent data. In this article, we take the first step to close this gap by introducing the FAMMA criterion and proving its asymptotic efficiency in the sense mentioned in the first paragraph. However, minimax optimality results are not pursued here. Our criterion has some distinctive features. First, it involves estimation of the high-dimensional inverse covariance matrix of a stationary time series that is not directly observable. Note that the covariance matrix of a stationary time series of length  $n$  can be viewed as a high-dimensional covariance matrix because its dimension is equivalent to the sample size. In situations where the error process is observable (or equivalently, the regression functions are known to be zero), Wu and Pourahmadi (2009) proposed a banded covariance matrix estimator of  $\Sigma_n$  and proved its consistency under spectral norm, which also leads to the consistency of the corresponding inverse matrix in estimating  $\Sigma_n^{-1}$ . These results were then extended by McMurry and Politis (2010) to tapered covariance matrix estimators. However, since the error process is in general unobservable, one can only estimate  $\Sigma_n^{-1}$  (or  $\Sigma_n$ ) through the output variables. As far as estimating  $\Sigma_n^{-1}$  is concerned, these output variables are contaminated by unknown regression functions. In Section 3, we propose estimating  $\Sigma_n^{-1}$  by its banded Cholesky decomposition with the corresponding parameters estimated *nonparametrically* from the *least squares residuals* of an increasing dimensional approximating model. We also obtain the rate of convergence of the proposed estimator, which plays a crucial role in proving the asymptotic efficiency of the FAMMA criterion. Second, our criterion is justified under a continuous weight set  $\mathcal{H}_N$  (see (2.10)). While  $\mathcal{H}_N$  is not as general as  $\mathcal{G}_n$ , as argued in Section 2, it can substantially reduce the limitations encountered by  $\mathcal{H}_n(N)$  when  $N$  is small.

It is worth mentioning that to justify MMA's asymptotic efficiency over the weight set  $\mathcal{G}_n$ , Wan, Zhang and Zou (2010) required a stringent condition on  $M$ ; see (2.20) of Section 2. As argued in Remark 4, this condition may preclude the approximating models whose estimators have the minimum risk (ignoring constants). When these models/estimators are precluded, the MMA criterion can only select weights for a set of suboptimal models/estimators, which is obviously not desirable. In fact, the same dilemma also arises in Ando and Li (2014), who used a similar assumption to prove their asymptotic efficiency results. Zhang, Wan and Zou (2013) considered model averaging problems in regression models with dependent errors. They adopted the JMA criterion to choose weights for a class of estimators and showed that the criterion is asymptotically efficient over the weight set  $\mathcal{G}_n$ . Their result, however, is still reliant on a condition similar to (2.20). In addition, the class of estimators considered in their paper, excluding all FGLS estimators, may suffer from lack of efficiency.

The remaining paper is organized as follows. In Section 2, we first concentrate on the case where  $\Sigma_n$  is known. We show in Theorem 1 that the autocovariance-corrected Mallows model averaging (AMMA) criterion, which is the FAMMA criterion with the estimator of  $\Sigma_n^{-1}$  replaced by  $\Sigma_n^{-1}$  itself, is asymptotically efficient. Since the assumptions used in Theorem 1 are rather mild, both Corollary 2.1 of Li (1987) and Theorem 1 of Hansen (2007) become its special case. We then turn attention to the more practical situation where  $\Sigma_n$  is unknown and propose choosing model weights by the FAMMA criterion. It is shown in Theorem 2 of Section 2 that the FAMMA criterion is asymptotically efficient as long as the corresponding estimator of  $\Sigma_n^{-1}$  has a sufficiently fast convergence rate. In Section 3, we provide a consistent estimator of  $\Sigma_n^{-1}$  based on its banded Cholesky decomposition, and derive the estimator's convergence rate under various situations. In Section 4, the asymptotic efficiency of the FAMMA criterion with  $\Sigma_n^{-1}$  estimated by the method

proposed in Section 3 is established. Finally, we conclude in Section 5. All proofs are relegated to the Appendix in order to maintain the flow of exposition.

## 2 The AMMA and FAMMA criteria

Consider a regression model with infinitely many parameters,

$$y_t = \sum_{j=1}^{\infty} \theta_j x_{tj} + e_t = \mu_t + e_t, t = 1, \dots, n, \quad (2.1)$$

where  $\mu_t = \sum_{j=1}^{\infty} \theta_j x_{tj}$ ,  $\mathbf{x}_t = (x_{t1}, x_{t2}, \dots)'$  is the explanatory vector with  $\sup_{t \geq 1, j \geq 1} E(x_{tj}^2) < \infty$ ,  $\theta_j, j \geq 1$  are unknown parameters satisfying  $\sum_{j=1}^{\infty} |\theta_j| < \infty$ , and  $\{e_t\}$ , independent of  $\{\mathbf{x}_t\}$ , is an unobservable stationary process with zero mean and finite variance. In matrix notation,  $Y_n = \boldsymbol{\mu}_n + \mathbf{e}_n$ , where  $Y_n = (y_1, \dots, y_n)'$ ,  $\boldsymbol{\mu}_n = (\mu_1, \dots, \mu_n)'$ , and  $\mathbf{e}_n = (e_1, \dots, e_n)'$ . The central focus of this paper is to explore how and to what extent the model averaging can be implemented in the presence of time series errors.

Let  $m = 1, \dots, M$  be a set of approximating models of (2.1), where the  $m$ th model uses the first  $k_m$  elements of  $\{\mathbf{x}_t\}$  with  $1 \leq k_1 < k_2 < \dots < k_M < n$  and  $M$  is allowed to grow to infinity with the sample size  $n$ . Assume that  $\Sigma_n = E(\mathbf{e}_n \mathbf{e}_n')$  is known and  $\Sigma_n^{-1}$  exists. Then the generalized least squares (GLS) estimator of the regression coefficient vector in the  $m$ th approximating model is given by  $\hat{\Theta}_m^* = (X_m' \Sigma_n^{-1} X_m)^{-1} X_m' \Sigma_n^{-1} Y_n$ , and the resultant estimate of  $\boldsymbol{\mu}_n$  is  $\hat{\boldsymbol{\mu}}_n(m) = P_m^* Y_n$ , where  $X_m = (x_{ij})_{1 \leq i \leq n, 1 \leq j \leq k_m}$ ,  $P_m^* = X_m (X_m' \Sigma_n^{-1} X_m)^{-1} X_m' \Sigma_n^{-1}$ , and  $X_M$  is assumed to be almost surely (a.s.) full rank throughout the paper. The model averaging estimator of  $\boldsymbol{\mu}_n$  based on the  $M$ th approximating models is  $\hat{\boldsymbol{\mu}}_n(\mathbf{w}) = P^*(\mathbf{w}) Y_n$ , where  $\mathbf{w} \in \mathcal{G}_n$  and  $P^*(\mathbf{w}) = \sum_{m=1}^M w_m P_m^*$ . To evaluate the performance of  $\hat{\boldsymbol{\mu}}_n(\mathbf{w})$ , we use the generalized squared error (GSE) loss

$$L_n^*(\mathbf{w}) = (\hat{\boldsymbol{\mu}}_n(\mathbf{w}) - \boldsymbol{\mu}_n)' \Sigma_n^{-1} (\hat{\boldsymbol{\mu}}_n(\mathbf{w}) - \boldsymbol{\mu}_n).$$

This loss function is a natural generalization of Hansen's (2007) average squared error in the sense that  $L_n^*(\mathbf{w})$  reduces to the latter when  $\Sigma_n^{-1}$  is replaced by the  $n \times n$  identity matrix. Through this generalization, it is easy to establish a connection between our results and some classical asymptotic efficiency results on model averaging/selection, thereby leading to a more comprehensive understanding of this research field. For further discussion, see Remarks 1 and 2 below. On the other hand, when the future values of  $y_t$  are entertained instead of the regression function  $\boldsymbol{\mu}_n$ , Wei and Yang (2012) proposed several different loss functions from a prediction point of view. Compared with the squared errors, their loss functions are particularly suitable for dealing with outliers.

The next lemma provides a representation for the conditional risk,

$$R_n^*(\mathbf{w}) = E(L_n^*(\mathbf{w}) | \mathbf{x}_1, \dots, \mathbf{x}_n) \equiv E_{\mathbf{x}}(L_n^*(\mathbf{w})),$$

which is an extension of Lemma 2 of Hansen (2007) to the case of dependent errors.

**Lemma 1.** Assume (2.1) and  $\Sigma_n^{-1}$  exists. Then, for any  $\mathbf{w} \in \mathcal{G}_n$ ,

$$R_n^*(\mathbf{w}) = \sum_{m=1}^M \sum_{l=1}^M w_m w_l [\boldsymbol{\mu}_n' \Sigma_n^{-1/2} (I - P_{\max\{m,l\}}) \Sigma_n^{-1/2} \boldsymbol{\mu}_n + \min\{k_m, k_l\}], \quad (2.2)$$

where  $P_j = \Sigma_n^{-1/2} X_j (X_j' \Sigma_n^{-1} X_j)^{-1} X_j' \Sigma_n^{-1/2}$  is the orthogonal projection matrix for the column space of  $\Sigma_n^{-1/2} X_j$ .

To choose a data-driven weight vector asymptotically minimizing  $L_n^*(\mathbf{w})$  over a suitable weight set  $\mathcal{H} \subseteq \mathcal{G}_n$ , we propose the AMMA criterion,

$$C_n^*(\mathbf{w}) = (Y_n - \hat{\boldsymbol{\mu}}_n(\mathbf{w}))' \Sigma_n^{-1} (Y_n - \hat{\boldsymbol{\mu}}_n(\mathbf{w})) + 2 \sum_{m=1}^M w_m k_m. \quad (2.3)$$

Note that the AMMA criterion is a special case of the criterion given in (6.2\*) of Andrews (1991) with  $M_n(h) = P^*(\mathbf{w})$  and  $W = \Sigma_n^{-1}$ . It reduces to the MMA criterion when

$$\{e_t\} \text{ is a sequence of i.i.d. random variables with } E(e_1) = 0, 0 < E(e_1^2) = \sigma^2 < \infty. \quad (2.4)$$

Recently, Liu, Okui and Yoshimura (2013) also suggested using  $C_n^*(\mathbf{w})$  to choose weight vectors in situations where  $e_t$  are independent but possibly heteroscedastic. While the AMMA criterion is not new to the literature, the question of whether its minimizer can (asymptotically) minimize  $L_n^*(\mathbf{w})$  seems rarely discussed, in particular when  $\mathcal{H}$  is uncountable.

Recall that by assuming (2.4),

$$\xi_n = \inf_{\mathbf{w} \in \mathcal{G}_n} R_n^*(\mathbf{w}) \rightarrow \infty \text{ a.s.}, \quad (2.5)$$

and

$$E(|e_1|^{4(N+1)} | \mathbf{x}_1) < \kappa < \infty \text{ a.s., for some positive integer } N, \quad (2.6)$$

Hansen (2007, Theorem 1) showed that  $C_n^*(\mathbf{w})$  is asymptotically efficient in the sense that

$$\frac{L_n^*(\bar{\mathbf{w}}_n)}{\inf_{\mathbf{w} \in \mathcal{H}_n(N)} L_n^*(\mathbf{w})} \rightarrow_p 1, \quad (2.7)$$

where  $\rightarrow_p$  denotes convergence in probability,

$$\mathcal{H}_n(N) = \{\mathbf{w} : w_m \in \{0, 1/N, 2/N, \dots, 1\}, \sum_{m=1}^M w_m = 1\}, \quad (2.8)$$

and

$$\bar{\mathbf{w}}_n = \arg \min_{\mathbf{w} \in \mathcal{H}_n(N)} C_n^*(\mathbf{w}).$$

Equation (2.7) gives a positive answer to the above question in the special case where  $\Sigma_n = \sigma^2 I_n$  and  $\mathcal{H} = \mathcal{H}_n(N)$  is a discrete set. When (2.6) holds for sufficiently large  $N$ , the restriction of  $\mathcal{G}_n$  to  $\mathcal{H}_n(N)$  is not an issue of overriding concern because the grid points  $i/N, i = 0, \dots, N$ , in  $\mathcal{H}_n(N)$  is dense enough to provide a good approximation for the optimal weight vector among

$$\bar{\mathcal{H}}_n(N) = \{\mathbf{w} : w_m \in [0, 1], 1 \leq \sum_{m=1}^M I_{\{w_m \neq 0\}} \leq N, \sum_{m=1}^M w_m = 1\}, \quad (2.9)$$

and hence among  $\mathcal{G}_n$ . We call  $\bar{\mathcal{H}}_n(N)$  *continuous extension* of  $\mathcal{H}_n(N)$  because it satisfies  $\mathcal{H}_n(N) \subseteq \bar{\mathcal{H}}_n(N)$  and  $a\mathbf{w}_1 + b\mathbf{w}_2 \in \bar{\mathcal{H}}_n(N)$ , for any  $0 \leq a, b \leq 1$  with  $a + b = 1$  and any  $\mathbf{w}_1 = (w_{11}, \dots, w_{M1})'$

and  $\mathbf{w}_2 = (w_{12}, \dots, w_{M2})' \in \bar{\mathcal{H}}_n(N)$  with  $\sum_{m=1}^M |I_{\{w_{m1} \neq 0\}} - I_{\{w_{m2} \neq 0\}}| = 0$ . It is shown in Hansen (2007, p.1179) that even when  $N = 2$ , the optimal weight vector in  $\bar{\mathcal{H}}_n(N)$  can yield an averaging estimator outperforming the one based on the optimal single model, except in some special cases.

On the other hand, when (2.6) holds only for moderate or small values of  $N$ , not only the optimal weight vector in  $\mathcal{G}_n$  but also that in  $\bar{\mathcal{H}}_n(N)$  cannot be well approximated by the elements in  $\mathcal{H}_n(N)$ . As a result, the advantage of model averaging over model selection becomes less apparent. To rectify this deficiency, we introduce the following continuous extension of  $\mathcal{H}_n(N)$ ,

$$\mathcal{H}_N = \bigcup_{l=1}^N \mathcal{H}_{(l)}, \quad (2.10)$$

where

$$\mathcal{H}_{(l)} = \{\mathbf{w} : \underline{\delta} \leq w_i I_{\{w_i \neq 0\}} \leq 1, \sum_{i=1}^M I_{\{w_i \neq 0\}} = l, \sum_{m=1}^M w_m = 1\},$$

with  $0 < \underline{\delta} < 1/N$ . The number of non-zero component is  $1 \leq l \leq N$  for any vector in  $\mathcal{H}_{(l)}$ . Therefore, this weight set leads to sparse combinations of  $\hat{\boldsymbol{\mu}}_n(m), m = 1, \dots, M$ . For a detailed discussion on sparse combinations from the minimax viewpoint, see Wang et al. (2014). We will show in Theorem 1 that

$$\frac{L_n^*(\tilde{\mathbf{w}}_n)}{\inf_{\mathbf{w} \in \mathcal{H}_N} L_n^*(\mathbf{w})} \rightarrow_p 1, \quad (2.11)$$

without the restriction  $\Sigma_n = \sigma^2 I_n$ , where

$$\tilde{\mathbf{w}}_n = \arg \inf_{\mathbf{w} \in \mathcal{H}_N} C_n^*(\mathbf{w}).$$

It is important to be aware that  $\mathcal{H}_N$  can inherit the benefits of  $\bar{\mathcal{H}}_n(N)$  mentioned previously because the difference between the two sets can be made arbitrarily small by making  $\underline{\delta}$  sufficiently close to 0. Technically speaking, a nonzero (regardless of how small)  $\underline{\delta}$  enables us to establish some sharp uniform probability bounds through replacing  $R_n^*(\mathbf{w})$  by suitable model selection risks (see (A.7)), thereby overcoming the difficulties arising from the uncountability of  $\mathcal{H}_N$ . In Remarks 3 and 4 after Theorem 1, we will also discuss the asymptotic efficiency of  $C_n^*(\mathbf{w})$  over more general weight sets such as  $\mathcal{G}_n$  and its variants. The following assumptions on  $\{e_t\}$  are needed in our analysis. As shown in the Appendix, these assumptions allow us to derive sharp bounds for the moments of quadratic forms in  $\{e_t\}$  using the first moment bound theorem of Findley and Wei (1993).

**Assumption 1.**  $\{e_t\}$  is a sequence of stationary time series with autocovariance function (ACF)  $\gamma_j = E(e_t e_{t+j})$  satisfying  $\sum_{j=-\infty}^{\infty} \gamma_j^2 < \infty$ , and admits a linear representation

$$e_t = \alpha_t + \sum_{k=1}^{\infty} \beta_k \alpha_{t-k} \quad (2.12)$$

in terms of the  $\mathcal{F}_t$ -measurable random variables  $\alpha_t$ , where  $\mathcal{F}_t, -\infty < t < \infty$  is an increasing sequence of  $\sigma$ -fields of events. Moreover,  $\{\alpha_t\}$  satisfies the following properties with probability 1:

$$(M1) \ E(\alpha_t | \mathcal{F}_{t-1}) = 0.$$

$$(M2) \ E(\alpha_t^2 | \mathcal{F}_{t-1}) = \sigma_\alpha^2.$$

(M3) There exist a positive integer  $N$  and a positive number  $S > 4N$  such that for some constant  $0 < C_S < \infty$ ,

$$\sup_{-\infty < t < \infty} E(|\alpha_t|^S | \mathcal{F}_{t-1}) \leq C_S. \quad (2.13)$$

**Assumption 2.** The spectral density function of  $\{e_t\}$ ,

$$f_e(\lambda) = (\sigma_\alpha^2/2\pi) \left| \sum_{j=0}^{\infty} \beta_j e^{-ij\lambda} \right|^2 \neq 0 \quad (2.14)$$

for all  $-\pi < \lambda \leq \pi$ , where  $\beta_0 = 1$ . Moreover,

$$\sum_{j=0}^{\infty} |\beta_j| < \infty. \quad (2.15)$$

We are now in a position to state Theorem 1.

**Theorem 1.** Assume Assumptions 1 and 2 in which  $N$  in (M3) is fixed. Let

$$D_n(m) = \boldsymbol{\mu}_n' \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \boldsymbol{\mu}_n + k_m \text{ and } k_n^* = \min_{1 \leq m \leq M} D_n(m).$$

Suppose

$$k_n^* \rightarrow \infty \text{ a.s.} \quad (2.16)$$

Then, (2.11) follows.

A few comments on Theorem 1 are in order.

**Remark 1.** Theorem 1 generalizes Theorem 1 of Hansen (2007) in several directions. First, (2.4) is a special case of (2.12), with  $\beta_k = 0$  for all  $k \geq 1$ . Second, the discrete weight set  $\mathcal{H}_n(N)$  is extended to its continuous extension  $\mathcal{H}_N$ . Third, when (2.4) holds and  $\{e_t\}$  is independent of  $\{\mathbf{x}_t\}$ , the moment condition (2.13) is milder than (2.6). Fourth, (2.5) is weakened to (2.16), which is much easier to verify. Note that  $k_n^*$  can be viewed as an index of the amount of information contained in the candidate models. Therefore, (2.16) is quite natural from the estimation theoretical viewpoint; see, e.g., Lai and Wei (1982), Yu, Lin and Cheng (2012) and Chan, Huang and Ing (2013). Suppose there exists a non-random and non-negative function  $Q(m)$  satisfying

$$\sup_{1 \leq m \leq M} \left| \frac{\boldsymbol{\mu}_n' \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \boldsymbol{\mu}_n}{n} - Q(m) \right| \rightarrow 0, \text{ a.s.}$$

Then, (2.16) is fulfilled if  $Q(m) \neq 0$  for all  $m$  and  $\lim_{m \rightarrow \infty} Q(m) = 0$ , which essentially require that all candidate models are misspecified, but those which have many parameters can give good approximations of the true model.



**Remark 2.** Corollary 2.1 of Li (1987) also becomes a special case of Theorem 1. To see this, note that under (2.4), (2.16) and

$$E(e_1^8) < \infty, \quad (2.17)$$

Li's (1987) Corollary 2.1 shows that (2.11) holds with  $N = 1$ , namely,  $C_n^*(\mathbf{w})$  is asymptotically efficient for model selection. However, since (2.4) and (2.17) also imply Assumptions 1 and 2, Li's conclusion readily follows from Theorem 1.

**Remark 3.** It is far from being trivial to extend (2.11) to

$$\frac{L_n^*(\mathbf{w}^0)}{\inf_{\mathbf{w} \in \mathcal{G}_n} L_n^*(\mathbf{w})} \rightarrow_p 1, \quad (2.18)$$

where

$$\mathbf{w}^0 = \arg \inf_{\mathbf{w} \in \mathcal{G}_n} C_n^*(\mathbf{w}).$$

Alternatively, if  $\alpha_t$  have light-tailed distributions, such as those described in (C2) and (C3) of Ing and Lai (2011), then by using the exponential probability inequalities developed in the same papers in place of the moment inequalities given in the proof of Theorem 1, it can be shown that (2.11) holds with  $\underline{\delta}$  tending to 0 and  $N = M$  tending to  $\infty$  sufficiently slowly with  $n$ . The details, however, are not reported here due to space constraints.

**Remark 4.** When (2.4) holds true, (2.18) has been developed in Theorem 1' of Wan, Zhang and Zou (2010) under

$$E(|e_t|^{4G} | \mathbf{x}_t) \leq \kappa < \infty \text{ a.s.}, \quad (2.19)$$

for some integer  $1 \leq G < \infty$ , and

$$M \xi_n^{-2G} \sum_{m=1}^M D_n^G(m) \rightarrow 0 \text{ a.s.} \quad (2.20)$$

Unfortunately, (2.20), imposing a stringent restriction on  $M$ , often precludes models having small GSE losses. To see this, assume that  $\mathbf{x}_t$  are nonrandom and the  $m$ th approximating model contains the first  $m$  regressors, namely  $k_m = m$ . Assume also that

$$D_n(m) = nm^{-a} + m \quad (2.21)$$

and the  $G$  in (2.19) is greater than  $1/a$  for some  $a \geq 1$ . It is easy to show that  $D_n(m)$  is minimized by  $m \sim (an)^{1/(1+a)}$ , yielding the optimal rate of  $D_n(m)$ ,  $n^{1/(1+a)}$ . Moreover, as will be clear from (A.5),  $D_n(m) = R_n^*(\mathbf{v}_m)$  is asymptotically equivalent to  $L_n^*(\mathbf{v}_m)$ , where  $\mathbf{v}_m$  is the  $m$ th standard unit vector in  $R^M$ . Hence the optimal rate of  $L_n^*(\mathbf{v}_m)$  is also  $n^{1/(1+a)}$ , which is achievable by any approximating model whose number of regressors  $m$  satisfying

$$c_1 n^{1/(1+a)} \leq m \leq c_2 n^{1/(1+a)}, \text{ for some } 0 < c_1 < c_2 < \infty. \quad (2.22)$$

If  $M \sim c_0 n^{1/(1+a)}$ , where  $c_0$  is any positive number, then for  $G > 1/a$  with  $a \geq 1$ , there exists  $c_3 > 0$  such that  $M \xi_n^{-2G} \sum_{m=1}^M D_n^G(m) \geq c_3 M n^G / n^{2G/(1+a)} \rightarrow \infty$  as  $n \rightarrow \infty$ , which violates (2.20).



In fact, it is shown in Example 2 of Wan, Zhang and Zou (2010) that a sufficient condition for (2.20) to hold is  $M = O(n^v)$  with  $v < G/(1 + 2aG) < 1/(1 + a)$ . These facts reveal that all models with  $L_n^*(\mathbf{v}_m)$  achieving the optimal rate  $n^{1/(1+a)}$  (or equivalently, with  $m$  obeying (2.22)) are excluded by (2.20). Under such a situation,  $C_n^*(\mathbf{w})$  can only select weights for a set of suboptimal models. Therefore, it is hard to conclude from their Theorem 1' that  $\hat{\boldsymbol{\mu}}_n(\mathbf{w}^0)$ 's GSE loss is asymptotically smaller than that of  $\hat{\boldsymbol{\mu}}_n(m)$  with  $m$  satisfying (2.22), even though this theorem guarantees  $C_n^*(\mathbf{w})$ 's asymptotic efficiency in the sense of (2.18). On the other hand, since Theorem 1 does not impose any restrictions similar to (2.20), one is free to choose  $M \sim \bar{C}n^{1/(1+a)}$ , with  $\bar{C}$  sufficiently large, so as to include the optimal model  $m \sim (an)^{1/(1+a)}$ . In addition, by noticing  $k_n^* = \min_{1 \leq m \leq M} D_n(m) \rightarrow \infty$  as  $n \rightarrow \infty$ , we know from Theorem 1 that  $\tilde{\mathbf{w}}_n$  satisfies (2.11), and hence  $\hat{\boldsymbol{\mu}}_n(\tilde{\mathbf{w}}_n)$  asymptotically outperforms the best one among  $\hat{\boldsymbol{\mu}}_n(m)$ ,  $1 \leq m \leq \bar{C}n^{1/(1+a)}$ , in terms of GSE loss. When  $a$  is unknown, the bound  $M \sim \bar{C}n^{1/(1+a)}$  is infeasible. However, if a strict lower bound for  $a$ , say  $\underline{a}$ , is known a priori, then the same conclusion still holds for  $M \sim \bar{C}^*n^{1/(1+\underline{a})}$  with any  $\bar{C}^* > 0$ .

**Remark 5.** It is worth noting that estimating the weight that minimizes  $L_n^*(\mathbf{w})$  will generally introduce a variance inflation factor, which may prevent us from obtaining the asymptotic efficiency. Under independent errors, a recent paper by Wang et al (2014) gives a comprehensive discussion of this matter from the minimax viewpoint. In fact, pursuing the minimax optimal rate is more relevant than the asymptotic efficiency in the presence of a large variance inflation factor. On the other hand, one can still attain the asymptotic efficiency by substantially suppressing this factor through: (i) reducing the size of the weight set and (ii) reducing the number of the candidate variables, which have been taken by Hansen (2007) and Wan et al. (2010), respectively. Unfortunately, the limitations imposed on the size of the weight set or  $M$  by these authors are too stringent, and hence may lead to suboptimal results, as discussed previously. Theorem 1 takes the first approach and provides a somewhat striking result that asymptotically efficient model averaging is still achievable under a continuous/uncountable weight set, which is in sharp contrast to Hansen's (2007) discrete/countable weight set. The theoretical underpinnings of Theorem 1 are some sharp uniform probability bounds, which are presented in the Appendix and established based on a mild lower bound condition on the weight set described in (2.10).

In the case where  $\Sigma_n$  is unknown, the asymptotic efficiency of  $C_n^*(\mathbf{w})$  developed in Theorem 1 becomes practically irrelevant. However, if there exists a consistent estimate,  $\hat{\Sigma}_n^{-1}$ , of  $\Sigma_n^{-1}$ , then the corresponding FGLS estimator of  $\boldsymbol{\mu}$  based on the  $m$ th approximating model is  $\hat{P}_m^* Y_n$ , where  $\hat{P}_m^* = X_m(X_m' \hat{\Sigma}_n^{-1} X_m)^{-1} X_m' \hat{\Sigma}_n^{-1}$ . Moreover, the FAMMA criterion,

$$\hat{C}_n^*(\mathbf{w}) = (Y_n - \hat{\boldsymbol{\mu}}_n^*(\mathbf{w}))' \hat{\Sigma}_n^{-1} (Y_n - \hat{\boldsymbol{\mu}}_n^*(\mathbf{w})) + 2 \sum_{m=1}^M w_m k_m, \quad (2.23)$$

can be used in place of  $C_n^*(\mathbf{w})$  to perform model averaging, where

$$\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) = \hat{P}^*(\mathbf{w}) Y_n = \sum_{m=1}^M w_m \hat{P}_m^* Y_n$$

is the FGLS model averaging estimator of  $\boldsymbol{\mu}_n$  given  $\mathbf{w}$ . Define

$$L_n^F(\mathbf{w}) = (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n)' \Sigma_n^{-1} (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n),$$

and

$$\hat{\mathbf{w}}_n = \arg \inf_{\mathbf{w} \in \mathcal{H}_N} \hat{C}_n^*(\mathbf{w}).$$

In the next theorem, we shall show that as long as  $\hat{\Sigma}_n^{-1}$  converges to  $\Sigma_n^{-1}$  sufficiently fast in terms of spectral norm, (2.11) still holds with  $L_n^*(\tilde{\mathbf{w}}_n)$  replaced by  $L_n^F(\hat{\mathbf{w}}_n)$ .

**Theorem 2.** *Assume that Assumptions 1 and 2 hold and there exists a sequence of positive numbers  $\{b_n\}$  satisfying  $b_n = o(n^{1/2})$  such that*

$$n\|\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}\|^2 = O_p(b_n^2), \quad (2.24)$$

where for the  $p \times p$  matrix  $A$ ,  $\|A\|^2 = \sup_{\mathbf{z} \in \mathbb{R}^p, \|\mathbf{z}\|=1} \mathbf{z}' A' A \mathbf{z}$  with  $\|\mathbf{z}\|$  denoting the Euclidean norm of  $\mathbf{z}$ . Moreover, suppose that there exists  $2N/S < \theta < 1/2$  such that

$$\lim_{n \rightarrow \infty} \frac{b_n^2}{k_n^{*1-2\theta}} = 0 \text{ a.s.} \quad (2.25)$$

Then,

$$\frac{L_n^F(\hat{\mathbf{w}}_n)}{\inf_{\mathbf{w} \in \mathcal{H}_N} L_n^*(\mathbf{w})} \rightarrow_p 1. \quad (2.26)$$

Below are some comments regarding Theorem 2.

**Remark 6.** In the next section, (2.24) will be established for  $\hat{\Sigma}_n^{-1} = \hat{\Sigma}_n^{-1}(q_n)$ , where  $\hat{\Sigma}_n^{-1}(q_n)$ , defined in (3.4), is obtained by the  $q_n$ -banded Cholesky decomposition of  $\Sigma_n^{-1}$  with the parameters in the Cholesky factors estimated nonparametrically from the least squares residuals of an increasing dimensional approximating model. As will be seen later, the order of the magnitude of  $b_n$  associated with  $\|\hat{\Sigma}_n^{-1}(q_n) - \Sigma_n^{-1}\|$  can vary depending on the strength of the dependence of  $\{e_t\}$ .

**Remark 7.** Zhang, Wan and Zou (2013) considered the model averaging estimator  $\tilde{\boldsymbol{\mu}}_n(\mathbf{w}) = \sum_{m=1}^M w_m \tilde{\boldsymbol{\mu}}_n(m)$  of  $\boldsymbol{\mu}$ , where

$$\tilde{\boldsymbol{\mu}}_n(m) = P_m^X Y_n \quad (2.27)$$

the estimator corresponding to the  $m$ th approximating model and  $P_m^X$ , an  $n \times n$  matrix, depends on  $\{\mathbf{x}_t\}$  only. They evaluated the performance of  $\tilde{\boldsymbol{\mu}}_n(\mathbf{w})$  using the usual squared error loss,

$$L_n(\mathbf{w}) = \|\tilde{\boldsymbol{\mu}}_n(\mathbf{w}) - \boldsymbol{\mu}\|^2,$$

and showed in their Theorem 2.1 that

$$\frac{L_n(\hat{\mathbf{w}}_n^{(J)})}{\inf_{\mathbf{w} \in \mathcal{G}_n} L_n(\mathbf{w})} \rightarrow_p 1, \quad (2.28)$$

where  $\hat{\mathbf{w}}_n^{(J)}$  is obtained from the JMA criterion (defined in equation (4) of their paper). While the weight set in (2.28) is more general than that in (2.26), an assumption similar to (2.20) is required in their proof of (2.28). In addition, (2.27), excluding all FGLS estimators (since  $\hat{\Sigma}_n^{-1}$  depends on

both  $\{\mathbf{x}_t\}$  and  $Y_n$ ), can suffer from lack of efficiency in estimating  $\boldsymbol{\mu}$ .

**Remark 8.** When  $e_t$  are independent random variables with  $E(e_t) = 0$  for all  $t$  and possibly unequal  $E(e_t^2) = \sigma_t^2$ , Liu, Okui and Yoshimura (2013, Theorem 4) obtained a weaker version of (2.26),

$$\frac{L_n^F(\hat{\mathbf{w}}_n)}{\inf_{\mathbf{w} \in \mathcal{H}_n(N)} L_n^*(\mathbf{w})} \rightarrow_p 1,$$

in which  $\hat{\Sigma}_n^{-1} = \text{diag}(\hat{\sigma}_1^{-2}, \dots, \hat{\sigma}_n^{-2})$  with  $\hat{\sigma}_t^{-2}$  satisfying

$$\sup_{1 \leq t \leq n} (\hat{\sigma}_t^{-2} - \sigma_t^{-2})^2 = O_p(n^{-1}), \quad (2.29)$$

among other conditions. However, since  $n^{-1}$  is a parametric rate, certain parametric assumptions on  $\sigma_t^2, 1 \leq t \leq n$ , are required to ensure (2.29). In addition, their proof, relying crucially on Theorem 2 of Whittle (1960), is not directly applicable to dependent data.

**Remark 9.** Assumption (2.25) is a strengthened version of (2.16). It essentially says that the (normalized) estimation error of  $\hat{\Sigma}_n^{-1}$  must be dominated by the amount of information contained in the candidate models in a certain way. This type of assumption seems indispensable for the FAMMA criterion to preserve the features of its infeasible counterpart.

**Remark 10.** Throughout this paper, the only assumption that we impose on  $\{\mathbf{x}_t\}$  is  $\sup_{t \geq 1, j \geq 1} E(|x_{tj}|^\nu) < \infty$  for some  $2 \leq \nu < \infty$ , in addition to the (a.s.) nonsingularity of  $X_M$ . Therefore,  $\{\mathbf{x}_t\}$  can be nonrandom, serially independent or serially dependent.

### 3 A consistent estimate of $\Sigma_n^{-1}$ based on the Cholesky decomposition.

In this section, we shall construct a consistent estimator of  $\Sigma_n^{-1}$  based on its banded Cholesky decomposition. Note first that according to (2.12), (2.14) and (2.15),  $e_t$  has an AR( $\infty$ ) representation,

$$\sum_{j=0}^{\infty} a_j e_{t-j} = \alpha_t, \quad (3.1)$$

where  $a_0 = 1$ ,  $\sum_{j=0}^{\infty} a_j z^j = (\sum_{j=0}^{\infty} \beta_j z^j)^{-1} \neq 0$  for all  $|z| \leq 1$  and  $\sum_{j=0}^{\infty} |a_j| < \infty$ ; see Zygmund (1959). If an AR( $k$ ),  $k \geq 1$ , model is used to approximate model (3.1), then the corresponding best (in the sense of mean squared error) AR coefficients are given by  $-(a_1(k), \dots, a_k(k))'$ , where

$$(a_1(k), \dots, a_k(k))' = \text{argmin}_{(c_1, \dots, c_k)' \in R^k} E(e_t + c_1 e_{t-1} + \dots + c_k e_{t-k})^2.$$

Define  $\sigma_k^2 = E(e_t + a_1(k)e_{t-1} + \dots + a_k(k)e_{t-k})^2$ . Then, the modified Cholesky decomposition for  $\Sigma_n^{-1}$  is

$$\Sigma_n^{-1} = \mathbf{T}_n' \mathbf{D}_n^{-1} \mathbf{T}_n, \quad (3.2)$$

where

$$\mathbf{D}_n = \text{diag}(\gamma_0, \sigma_1^2, \sigma_2^2, \dots, \sigma_{n-1}^2),$$

and  $\mathbf{T}_n = (t_{ij})_{1 \leq i, j \leq n}$  is a lower triangular matrix satisfying

$$t_{ij} = \begin{cases} 0, & \text{if } i < j; \\ 1, & \text{if } i = j; \\ a_{i-j}(i-1), & \text{if } 2 \leq i \leq n, 1 \leq j \leq i-1. \end{cases}$$

Since  $\mathbf{T}_n$  and  $\mathbf{D}_n$  may contain too many parameters as compared with  $n$ , we are led to consider a banded Cholesky decomposition of  $\Sigma_n^{-1}$ ,

$$\Sigma_n^{-1}(q) = \mathbf{T}'_n(q) \mathbf{D}_n^{-1}(q) \mathbf{T}_n(q), \quad (3.3)$$

where  $1 \leq q \ll n$  is referred to as the banding parameter,

$$\mathbf{D}_n(q) = \text{diag}(\gamma_0, \sigma_1^2, \dots, \sigma_q^2, \dots, \sigma_q^2),$$

and  $\mathbf{T}_n(q) = (t_{ij}(q))_{1 \leq i, j \leq n}$  with

$$t_{ij}(q) = \begin{cases} 0, & \text{if } i < j \text{ or } \{q+1 < i \leq n, 1 \leq j \leq i-q-1\}; \\ 1, & \text{if } i = j; \\ a_{i-j}(i-1), & \text{if } 2 \leq i \leq q, 1 \leq j \leq i-1; \\ a_{i-j}(q), & \text{if } q+1 \leq i \leq n, i-q \leq j \leq i-1. \end{cases}$$

To estimate the banded Cholesky factors in (3.3), we first generate the least squares residuals  $\hat{\mathbf{e}}_n = (\hat{e}_1, \dots, \hat{e}_n)'$  based on the approximating model  $\sum_{j=1}^d \theta_j x_{t,j}$  for (2.1), where  $d = d_n$  is allowed to grow to infinity with  $n$  and  $\hat{\mathbf{e}}_n = (I - H_d)Y_n$  with  $H_d$  denoting the orthogonal projection matrix for the column space of  $X(d) = (x_{ij})_{1 \leq i \leq n, 1 \leq j \leq d}$ . Having obtained  $\hat{\mathbf{e}}_n$ , the  $\gamma_0$  and  $\sigma^2(k)$  in  $\mathbf{D}_n(q)$  and the  $a_i(k)$  in  $\mathbf{T}_n(q)$  can be estimated by  $\hat{\gamma}_0$ ,  $\hat{\sigma}^2(k)$ , and  $\hat{a}_i(k)$ , respectively, where for  $1 \leq k \leq q$ ,

$$\begin{aligned} (\hat{a}_1(k), \dots, \hat{a}_k(k))' &= \underset{(c_1, \dots, c_k)' \in R^k}{\text{argmin}} \sum_{t=q+1}^n (\hat{e}_t + c_1 \hat{e}_{t-1} + \dots + c_k \hat{e}_{t-k})^2, \\ \hat{\gamma}_0 &= n^{-1} \sum_{t=1}^n \hat{e}_t^2, \\ \hat{\sigma}_k^2 &= (n-q)^{-1} \sum_{t=q+1}^n (\hat{e}_t + \sum_{j=1}^k \hat{a}_j(k) \hat{e}_{t-j})^2. \end{aligned}$$

Plugging these estimators into  $\mathbf{D}_n(q)$  and  $\mathbf{T}_n(q)$ , we obtain  $\hat{\mathbf{D}}_n(q)$  and  $\hat{\mathbf{T}}_n(q)$ , and hence an estimator of  $\Sigma_n^{-1}$ ,

$$\hat{\Sigma}_n^{-1}(q) = \hat{\mathbf{T}}'_n(q) \hat{\mathbf{D}}_n^{-1}(q) \hat{\mathbf{T}}_n(q). \quad (3.4)$$

Note that we have suppressed the dependence of  $\hat{\Sigma}_n^{-1}(q)$  on  $d$  in order to simplify notation. The next theorem provides a rate of convergence of  $\hat{\Sigma}_n^{-1}(q)$  to  $\Sigma_n^{-1}$  when  $q = q_n$  and  $d = d_n$  grow to infinity with  $n$  at suitable rates.

**Theorem 3.** Assume Assumptions 1 and 2 with (2.13) and (2.15) replaced by

$$\mathbb{E}(|\alpha_t|^{2r}|\mathcal{F}_{t-1}) < C_{2r} < \infty \text{ a.s.}, \quad (3.5)$$

where  $r \geq 2$ , and

$$\sum_{j \geq 1} j|\beta_j| < \infty, \quad (3.6)$$

respectively. Also assume that

$$\sup_{t \geq 1, j \geq 1} \mathbb{E}(|x_{tj}|^{2r}) < \infty. \quad (3.7)$$

Suppose that  $d_n$  and  $q_n$  are chosen to satisfy:

$$d_n \asymp n^{1/4}, \quad (3.8)$$

$$\max\{d_n, q_n\} \sum_{j \geq d_n} |\theta_j| = o(1), \quad (3.9)$$

and

$$\frac{q_n^2 d_n}{n} = o(1). \quad (3.10)$$

Then,

$$\|\hat{\Sigma}_n^{-1}(q_n) - \Sigma_n^{-1}\| = O_p \left( \frac{q_n^{1+r-1}}{n^{1/2}} + \sqrt{\sum_{j \geq q_n+1} |a_j| \sum_{j \geq q_n+1} j|a_j|} \right). \quad (3.11)$$

**Remark 11.** In the simpler situation where  $\hat{\mathbf{e}}_n = Y_n = \mathbf{e}_n$ , namely,  $\mu_t = 0$  for all  $t$ , Wu and Pourahmadi (2009) proposed a banded covariance matrix estimator  $\check{\Sigma}_{n,l} = (\hat{\gamma}_{i-j} I_{|i-j| \leq l})_{1 \leq i, j \leq n}$  of  $\Sigma_n$ , where  $\hat{\gamma}_k = n^{-1} \sum_{i=1}^{n-|k|} e_i e_{i+|k|}$  is the  $k$ th lag sample ACF of  $\{e_t\}$  and  $l$  is also called the banding parameter. When  $l = l_n = o(n^{1/2})$  and (3.5) holds with  $r = 2$ , their Theorems 2 and 3 imply that  $\check{\Sigma}_{n,l_n}$  is positive definite with probability approaching one,

$$\|\check{\Sigma}_{n,l_n} - \Sigma_n\| = O_p \left( \frac{l_n}{n^{1/2}} + \sum_{j \geq q_n} |\gamma_j| \right), \quad (3.12)$$

and

$$\|\check{\Sigma}_{n,l_n}^{-1} - \Sigma_n^{-1}\| = O_p \left( \frac{l_n}{n^{1/2}} + \sum_{j \geq q_n} |\gamma_j| \right). \quad (3.13)$$

McMurry and Politis (2010) generalized (3.12) and (3.13) to tapered covariance matrix estimators. Ing, Chiou and Guo (2013) considered estimating  $\Sigma_n^{-1}$  through the banded Cholesky decomposition approach in situations where  $\hat{\mathbf{e}}_n$  is obtained by a *correctly specified* regression model. They

established the consistency of the proposed estimator under spectral norm, even when  $\{e_t\}$  is a long-memory time series. However, since this section allows the regression model to be misspecified, all the aforementioned results are not directly applicable here.

**Remark 12.** The second term on the right-hand side of (3.11) is mainly contributed by the approximation error  $\|\Sigma_n^{-1}(q_n) - \Sigma_n^{-1}\|$ , whereas the first one is mainly due to the sampling variability  $\|\hat{\Sigma}_n^{-1}(q_n) - \Sigma_n^{-1}(q_n)\|$ , which is in turn dominated by  $\|\hat{\mathbf{T}}_n(q_n) - \mathbf{T}_n(q_n)\|$ , as shown in the proof of Theorem 3. Similarly, the first and second terms on the right-hand side of (3.12) are contributed by  $\|\check{\Sigma}_{n,l_n} - \Sigma_{n,l_n}\|$  and  $\|\Sigma_{n,l_n} - \Sigma_n\|$ , respectively. Here,  $\Sigma_{n,l_n} = (\gamma_{i-j} I_{|i-j| \leq l_n})_{1 \leq i,j \leq n}$  is the population version of  $\check{\Sigma}_{n,l_n}$ . However, unlike  $\check{\Sigma}_{n,l_n} - \Sigma_{n,l_n}$ ,  $\hat{\mathbf{T}}_n(q_n) - \mathbf{T}_n(q_n)$  is not a Toeplitz matrix. Hence our upper bound for  $\|\hat{\mathbf{T}}_n(q_n) - \mathbf{T}_n(q_n)\|$  is derived from complicated maximal probability inequalities, such as (A.42) and (A.49), which also lead to an additional exponent  $r^{-1}$  in the first term on the right-hand side of (3.11).

**Remark 13.** The technical assumptions (3.8)-(3.10) essentially say that the dimension,  $d_n$ , of the working regression model shouldn't be too large or too small. They ensure that the sampling variability and the approximation error introduced by this model are completely absorbed into the first or second term on the right-hand side of (3.11), which depend only on the working AR model used in the Cholesky decomposition. As shown in the next section, this feature can substantially reduce the burden of verifying (2.24) and (2.25).

## 4 Asymptotic efficiency of the FAMMA method with $\hat{\Sigma}_n^{-1} = \hat{\Sigma}_n^{-1}(q_n)$ .

In this section, we shall establish the asymptotic efficiency of  $\hat{C}_n^*(\mathbf{w})$  with  $\hat{\Sigma}_n^{-1} = \hat{\Sigma}_n^{-1}(q_n)$ , denoted by  $\hat{C}_{n,q_n}^*(\mathbf{w})$ , when the AR coefficients of  $\{e_t\}$  satisfy

$$\sqrt{\sum_{j \geq q} |a_j| \sum_{j \geq q} j |a_j|} \leq C_1 \exp(-\nu q), \quad (4.1)$$

or

$$\sqrt{\sum_{j \geq q} |a_j| \sum_{j \geq q} j |a_j|} \leq C_2 q^{-\nu}, \quad (4.2)$$

for all  $q \geq 1$  and some positive constants  $C_1, C_2$  and  $\nu$ . We call (4.1) the exponential decay case, which is fulfilled by any causal and invertible ARMA( $p, q$ ) model with  $0 \leq p, q < \infty$ . On the other hand, (4.2) is referred to as the algebraic decay case, which is commonly discussed in the context of model selection for time series; see Shibata (1981) and Ing and Wei (2003, 2005).

We first choose suitable  $q_n$  for  $\hat{\Sigma}_n^{-1}(q_n)$  to ensure that the bound in (3.11) possesses the optimal rate. When (4.1) is assumed, it is not difficult to see that the optimal rate of (3.11) is  $O_p((\log n)^{1+r^{-1}}/n^{1/2})$ , which is achieved by

$$q_n = c_4 \log n, \quad (4.3)$$

for some sufficiently large constant  $c_4$ . Therefore, (2.24) holds with

$$\hat{\Sigma}_n^{-1} = \hat{\Sigma}_n^{-1}(c_4 \log n) \text{ and } b_n = (\log n)^{1+r^{-1}}. \quad (4.4)$$

When (4.2) is true, by letting

$$q_n = \lfloor n^{1/\{2(1+r^{-1}+\nu)\}} \rfloor, \quad (4.5)$$

where  $\lfloor a \rfloor$  denotes the largest integer  $\leq a$ , we get the optimal rate of (3.11),  $O_p(n^{-\nu/\{2(1+\nu+r^{-1})\}})$ , yielding that (2.24) holds with

$$\hat{\Sigma}_n^{-1} = \hat{\Sigma}_n^{-1}(\lfloor n^{1/\{2(1+r^{-1}+\nu)\}} \rfloor) \text{ and } b_n = n^{\frac{1+r^{-1}}{2(1+\nu+r^{-1})}}. \quad (4.6)$$

We are ready to establish the asymptotic efficiency of  $\hat{C}_{n,q_n}^*(\mathbf{w})$  under (4.1).

**Corollary 1.** *Assume Assumptions 1 and 2, (4.1) and (3.7) with  $2r$  replaced by  $S$ , noting that  $S$  is defined in (M3) of Assumption 1. Suppose that  $d_n$  and  $q_n$  obey (3.8) and (4.3), respectively. Moreover, assume*

$$d_n \sum_{j \geq d_n} |\theta_j| = o(1), \quad (4.7)$$

and for some  $2N/S < \theta < 1/2$ ,

$$\frac{(\log n)^{1+(2/S)}}{k_n^{*(1/2)-\theta}} = 0 \text{ a.s.} \quad (4.8)$$

Then, (2.26) holds with  $\hat{\mathbf{w}}_n = \arg \inf_{\mathbf{w} \in \mathcal{H}_N} \hat{C}_{n,q_n}^*(\mathbf{w})$ .

Corollary 1 follows directly from Theorems 2 and 3 and (4.4) with  $r^{-1}$  replaced by  $2/S$ . Its proof is thus omitted. Condition (4.8) is easily satisfied when  $D_n(m)$  follows (2.21). To see this, note that (2.21) implies  $k_n^* = c_5 n^{1/(1+a)}$  for some  $c_5 > 0$ . Therefore, (4.8) holds for any  $2N/S < \theta < 1/2$ . On the other hand, it is not difficult to show that (4.8) is violated when  $D(m) = n \exp(-c_6 m) + m$  for some  $c_6 > 0$ , which leads to a much smaller  $k_n^* = c_7 \log n$  for some  $c_7 > 0$ .

To establish the asymptotic efficiency of  $\hat{C}_{n,q_n}^*(\mathbf{w})$  under (4.2) with  $\nu$  unknown, we need to assume that  $\nu$  has a known lower limit  $\nu_0 \geq 1/3$ .

**Corollary 2.** *Assume Assumptions 1 and 2, (4.2) and (3.7) with  $2r$  replaced by  $S$ . Suppose that  $d_n$  obeys (3.8) and  $q_n$  satisfies (4.5) with  $r^{-1}$  replaced by  $2/S$  and  $\nu$  by  $\nu_0$ . Moreover, assume (3.9) and for some  $2N/S < \theta < 1/2$ ,*

$$\frac{n^{\frac{1+(2/S)}{2[1+\nu_0+(2/S)]}}}{k_n^{*(1/2)-\theta}} = 0 \text{ a.s.} \quad (4.9)$$

Then, the conclusion of Corollary 1 follows.

Corollary 2 can be proved using Theorems 2 and 3 and (4.6) with  $r^{-1}$  and  $\nu$  replaced by  $2/S$  and  $\nu_0$ , respectively. We again omit the details. Before closing this section, we provide a sufficient condition for (4.9) in situations where  $D_n(m)$  obeys (2.21). We assume that (2.13) in Assumption 1 holds for any  $0 < S < \infty$  in order to simplify exposition. Elementary calculations show that (4.9) follows from  $\nu_0 > a$ . However, since  $a$  is in general unknown, our simple and practical guidance for verifying (4.9) is to check whether  $\nu_0 > \bar{a}$ , where  $\bar{a}$  is a known upper bound for  $a$ .



## 5 Concluding remarks

This paper provides guidance for the model averaging implementation in regression models with time series errors. Driven by the efficiency improvement, our goal is to choose the optimal weight vector that averages across FGLS estimators obtained from a set of approximating models of the true regression function. We propose the FAMMA as the weight selection criterion and show its asymptotic optimality in the sense of (2.26). To the best of our knowledge, it is the first time that the FGLS-based criterion is proved to have this type of property in the presence of time-dependent errors.

On the other hand, our asymptotic optimality, implicitly involving the search for the averaging estimator whose loss (or conditional risk) has the best constant in addition to the best rate, is typically not achievable when the number of candidate models is large and the models are not necessarily nested. While Wan, Zhang and Zou (2010) and Zhang, Wan and Zou (2013) proved the asymptotic efficiency of their averaging estimators without assuming nested candidate models, a stringent condition on the number of models, e.g., (2.20), is placed as the tradeoff. Furthermore, on top of their positive report, no clear guideline for the optimal averaging across arbitrary combinations of regressors was offered. In fact, in this more challenging situation, pursuing the minimax optimal rate appears to be more relevant than the asymptotic efficiency. The theoretical results developed in Wang et al. (2014) and in Sections 2 and 3 provide useful tools for deriving the minimax optimal rate under model (2.1). Moreover, motivated by Ing and Lai (2011), we conjecture that when the variables are preordered by the orthogonal greedy algorithm (OGA) (see, e.g., Temlyakov (2000) and Ing and Lai (2011)), this rate is achievable by FAMMA with 2 replaced by a factor directly proportional to the natural logarithm of the number of candidate models. We leave investigations along this research direction to future work.

## APPENDIX

*Proof of Lemma 1.* Note first that  $R_n^*(\mathbf{w}) = E_{\mathbf{x}}(L_n^*(\mathbf{w})) = E_{\mathbf{x}}(\mathbf{e}_n' P^*(\mathbf{w}) \Sigma_n^{-1} P^*(\mathbf{w}) \mathbf{e}_n) + \boldsymbol{\mu}_n' (I - P^*(\mathbf{w}))' \Sigma_n^{-1} (I - P^*(\mathbf{w})) \boldsymbol{\mu}_n$ . Since

$$\Sigma_n^{-1/2} P^*(\mathbf{w}) = \sum_{m=1}^M w_m P_m \Sigma_n^{-1/2}, \quad (\text{A.1})$$

it follows that

$$\begin{aligned} & E_{\mathbf{x}}(\mathbf{e}_n' P^*(\mathbf{w}) \Sigma_n^{-1} P^*(\mathbf{w}) \mathbf{e}_n) \\ &= E_{\mathbf{x}}\left(\sum_{m=1}^M \sum_{l=1}^M w_m w_l \mathbf{e}_n' \Sigma_n^{-1/2} P_l P_m \Sigma_n^{-1/2} \mathbf{e}_n\right) \\ &= \sum_{m=1}^M \sum_{l=1}^M w_l w_m \min\{k_m, k_l\}. \end{aligned}$$

Similarly,  $\boldsymbol{\mu}_n' (I - P^*(\mathbf{w}))' \Sigma_n^{-1} (I - P^*(\mathbf{w})) \boldsymbol{\mu}_n = \sum_{m=1}^M \sum_{l=1}^M w_l w_m \boldsymbol{\mu}_n' \Sigma_n^{-1/2} (I - P_{\max\{m,l\}}) \Sigma_n^{-1/2} \boldsymbol{\mu}_n$ . Consequently, the desired conclusion (2.2) follows.

*Proof of Theorem 1.* Define  $\mathbf{w}_n^* = \arg \min_{\mathbf{w} \in \mathcal{H}_N} L_n^*(\mathbf{w})$ ,  $(\tilde{w}_{n,1}, \dots, \tilde{w}_{n,M})' = \tilde{\mathbf{w}}_n$ , and  $(w_{n,1}^*, \dots, w_{n,M}^*)' =$

$\mathbf{w}_n^*$ . By noticing

$$\begin{aligned} C_n^*(\mathbf{w}) - L_n^*(\mathbf{w}) &= \mathbf{e}_n' \Sigma_n^{-1} \mathbf{e}_n + 2\mathbf{e}_n' \Sigma_n^{-1} (I - P^*(\mathbf{w})) \boldsymbol{\mu}_n \\ &- 2\{\mathbf{e}_n' \Sigma_n^{-1} P^*(\mathbf{w}) \mathbf{e}_n - \sum_{m=1}^M w_m k_m\}, \end{aligned}$$

we get

$$\begin{aligned} 0 &\geq \{C_n^*(\tilde{\mathbf{w}}_n) - C_n^*(\mathbf{w}_n^*)\} = L_n^*(\tilde{\mathbf{w}}_n) - L_n^*(\mathbf{w}_n^*) + 2\mathbf{e}_n' \Sigma_n^{-1} (I - P^*(\tilde{\mathbf{w}}_n)) \boldsymbol{\mu}_n \\ &- 2\left\{\mathbf{e}_n' \Sigma_n^{-1} P^*(\tilde{\mathbf{w}}_n) \mathbf{e}_n - \sum_{m=1}^M \tilde{w}_{n,m} k_m\right\} - 2\mathbf{e}_n' \Sigma_n^{-1} (I - P^*(\mathbf{w}_n^*)) \boldsymbol{\mu}_n \\ &+ 2\left\{\mathbf{e}_n' \Sigma_n^{-1} P^*(\mathbf{w}_n^*) \mathbf{e}_n - \sum_{m=1}^M w_{n,m}^* k_m\right\} \\ &= L_n^*(\tilde{\mathbf{w}}_n) - L_n^*(\mathbf{w}_n^*) + 2A_n(\tilde{\mathbf{w}}_n) - 2B_n(\tilde{\mathbf{w}}_n) - 2A_n(\mathbf{w}_n^*) + 2B_n(\mathbf{w}_n^*), \end{aligned} \quad (\text{A.2})$$

where  $A_n(\mathbf{w}) = \mathbf{e}_n' \Sigma_n^{-1} (I - P^*(\mathbf{w})) \boldsymbol{\mu}_n$  and  $B_n(\mathbf{w}) = \mathbf{e}_n' \Sigma_n^{-1} P^*(\mathbf{w}) \mathbf{e}_n - \sum_{m=1}^M w_m k_m$ . In view of (A.2) and  $L_n^*(\tilde{\mathbf{w}}_n) \geq L_n^*(\mathbf{w}_n^*)$ , it suffices for (2.11) to show that

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{A_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1), \quad (\text{A.3})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{B_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1), \quad (\text{A.4})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{L_n^*(\mathbf{w})}{R_n^*(\mathbf{w})} - 1 \right| = o_p(1), \quad (\text{A.5})$$

where  $\xrightarrow{p}$  denotes convergence in probability.

To show (A.3), first note that

$$\mathcal{H}_{(l)} = \bigcup_{1 \leq j_1 < j_2 < \dots < j_l \leq M} \mathcal{H}_{j_1, \dots, j_l},$$

where for  $1 \leq j_1 < \dots < j_l \leq M$ ,  $\mathcal{H}_{j_1, \dots, j_l} = \{\mathbf{w} : \mathbf{w} \in \mathcal{H}_{(l)} \text{ and } \omega_{j_i} \neq 0, 1 \leq i \leq l\}$ . Hence for any  $\varepsilon > 0$ ,

$$\begin{aligned} &P_{\mathbf{x}} \left( \sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{A_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| > \varepsilon \right) \\ &\leq \sum_{l=1}^N \sum_{j_l=l}^M \dots \sum_{j_1=1}^{j_2-1} P_{\mathbf{x}} \left( \sup_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} \left| \frac{A_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| > \varepsilon \right) \\ &\leq \sum_{l=1}^N \sum_{j_l=l}^M \dots \sum_{j_1=1}^{j_2-1} P_{\mathbf{x}} \left( \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}_n' \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \mathbf{e}_n \right|}{\underline{\delta}^2 \max_{m \in \{j_1, \dots, j_l\}} D_n(m)} > \varepsilon \right) \equiv \sum_{l=1}^N Q_l, \end{aligned} \quad (\text{A.6})$$

where  $P_{\mathbf{x}}(\cdot) = P(\cdot | \mathbf{x}_1, \dots, \mathbf{x}_n)$  and the second inequality follows from

$$\inf_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} R_n^*(\mathbf{w}) \geq \underline{\delta}^2 \max_{m \in \{j_1, \dots, j_l\}} D_n(m), \quad (\text{A.7})$$

which is ensured by Lemma 1 and the definition of  $\mathcal{H}_N$ . Let  $S_1 = S/2$ . Then, by Chebyshev's inequality, (M3), and Lemma 2 of Wei (1987), it holds that

$$\begin{aligned} & P_{\mathbf{x}} \left( \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \mathbf{e}_n \right|}{\underline{\delta}^2 \max_{m \in \{j_1, \dots, j_l\}} D_n(m)} > \varepsilon \right) \\ & \leq C \sum_{m \in \{j_1, \dots, j_l\}} \left( \frac{\mathbb{E}_{\mathbf{x}} \{ \boldsymbol{\mu}'_n \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \mathbf{e}_n \}^{S_1}}{\left\{ \max_{m \in \{j_1, \dots, j_l\}} D_n(m) \right\}^{S_1}} \right) \\ & \leq C \sum_{m \in \{j_1, \dots, j_l\}} \frac{\left( \boldsymbol{\mu}'_n \Sigma_n^{-1/2} (I - P_m) \Sigma_n^{-1/2} \boldsymbol{\mu}_n \right)^{S_1/2}}{\left( \max_{m \in \{j_1, \dots, j_l\}} D_n(m) \right)^{S_1}} \\ & \leq C \sum_{m \in \{j_1, \dots, j_l\}} \frac{1}{\left( \max_{m \in \{j_1, \dots, j_l\}} D_n(m) \right)^{S_1/2}} \leq C \frac{l}{(D_n(j_l))^{S_1/2}}, \end{aligned}$$

where here and hereafter  $C$  denotes a generic positive constant whose value is independent of  $n$  and may vary at different occurrences. Therefore, for each  $1 \leq l \leq N$ ,

$$\begin{aligned} Q_l & \leq C \left\{ \sum_{j_l=l}^{k_n^*} \cdots \sum_{j_1=1}^{j_2-1} \frac{1}{(k_n^*)^{S_1/2}} + \sum_{j_l=k_n^*+1}^M \cdots \sum_{j_1=1}^{j_2-1} \frac{1}{(D_n(j_l))^{S_1/2}} \right\} \\ & \leq C \left\{ k_n^{*(S_1/2-l)} + \sum_{j_l=k_n^*+1}^{\infty} \frac{j_l^{l-1}}{j_l^{S_1/2}} \right\}, \end{aligned}$$

which converges to 0 a.s. in view of (2.16). As a result,

$$P_{\mathbf{x}} \left( \sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{A_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| > \varepsilon \right) \rightarrow 0, \text{ a.s.}$$

This and the dominated convergence theorem together imply (A.3).

Similarly,

$$P_{\mathbf{x}} \left( \sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{B_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| > \varepsilon \right) \leq C \sum_{l=1}^N E_l,$$

where

$$E_l = \sum_{j_l=l}^M \cdots \sum_{j_1=1}^{j_2-1} P_{\mathbf{x}} \left( \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left| \mathbf{e}'_n \Sigma_n^{-1/2} P_m \Sigma_n^{-1/2} \mathbf{e}_n - k_m \right|}{\underline{\delta}^2 \max_{m \in \{j_1, \dots, j_l\}} D_n(m)} > \varepsilon \right).$$

By (M3) and the first moment bound theorem of Findley and Wei (1993), it follows that

$$\begin{aligned} & P_{\mathbf{x}} \left( \frac{\sum_{m \in \{j_1, \dots, j_l\}} |\mathbf{e}'_n \Sigma_n^{-1/2} P_m \Sigma_n^{-1/2} \mathbf{e}_n - k_m|}{\delta^2 \max_{m \in \{j_1, \dots, j_l\}} D_n(m)} > \varepsilon \right) \\ & \leq C \sum_{m \in \{j_1, \dots, j_l\}} \frac{k_m^{S_1/2}}{\left( \max_{m \in \{j_1, \dots, j_l\}} D_n(m) \right)^{S_1}} \leq \frac{Cl}{(D_n(j_l))^{S_1/2}}. \end{aligned}$$

Therefore, (A.4) follows immediately from an argument similar to that used to prove (A.3). The proof of (A.5) is similar to those of (A.3) and (A.4). The details are omitted.

*Proof of Theorem 2.* Define  $\hat{L}_n^*(\mathbf{w}) = (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n)' \hat{\Sigma}_n^{-1} (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n)$ . Then, it follows that

$$\begin{aligned} \hat{C}_n^*(\mathbf{w}) &= (\mathbf{e}_n - (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n))' \hat{\Sigma}_n^{-1} (\mathbf{e}_n - (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n)) + 2 \sum_{m=1}^M w_m k_m \\ &= \mathbf{e}'_n \hat{\Sigma}_n^{-1} \mathbf{e}_n + \hat{L}_n^*(\mathbf{w}) - 2 (\hat{\boldsymbol{\mu}}_n^*(\mathbf{w}) - \boldsymbol{\mu}_n)' \hat{\Sigma}_n^{-1} \mathbf{e}_n + 2 \sum_{m=1}^M w_m k_m \\ &= \mathbf{e}'_n \hat{\Sigma}_n^{-1} \mathbf{e}_n + 2 \boldsymbol{\mu}'_n (I - \hat{P}^*(\mathbf{w}))' \hat{\Sigma}_n^{-1} \mathbf{e}_n + \hat{L}_n^*(\mathbf{w}) - 2 \left\{ \mathbf{e}'_n \hat{P}^{*'}(\mathbf{w}) \hat{\Sigma}_n^{-1} \mathbf{e}_n - \sum_{m=1}^M w_m k_m \right\}, \end{aligned}$$

and hence

$$\begin{aligned} 0 &\geq \hat{C}_n^*(\hat{\mathbf{w}}_n) - \hat{C}_n^*(\mathbf{w}_n^F) \\ &= \hat{L}_n^*(\hat{\mathbf{w}}_n) - \hat{L}_n^*(\mathbf{w}_n^F) + 2 \hat{A}_n(\hat{\mathbf{w}}_n) - 2 \hat{A}_n(\mathbf{w}_n^F) - 2 \hat{B}_n(\hat{\mathbf{w}}_n) + 2 \hat{B}_n(\mathbf{w}_n^F) \\ &= (\hat{L}_n^*(\hat{\mathbf{w}}_n) - L_n^F(\hat{\mathbf{w}}_n)) - (\hat{L}_n^*(\mathbf{w}_n^F) - L_n^F(\mathbf{w}_n^F)) + 2 \hat{A}_n(\hat{\mathbf{w}}_n) - 2 \hat{A}_n(\mathbf{w}_n^F) \\ &\quad - 2 \hat{B}_n(\hat{\mathbf{w}}_n) + 2 \hat{B}_n(\mathbf{w}_n^F) + (L_n^F(\hat{\mathbf{w}}_n) - L_n^F(\mathbf{w}_n^F)), \end{aligned}$$

where  $\mathbf{w}_n^F = \arg \min_{\mathbf{w} \in \mathcal{H}_N} L_n^F(\mathbf{w})$ ,  $\hat{A}_n(\mathbf{w}) = \boldsymbol{\mu}'_n (I - \hat{P}^*(\mathbf{w}))' \hat{\Sigma}_n^{-1} \mathbf{e}_n$  and  $\hat{B}_n(\mathbf{w}) = \mathbf{e}'_n \hat{P}^{*'}(\mathbf{w}) \hat{\Sigma}_n^{-1} \mathbf{e}_n - \sum_{m=1}^M w_m k_m$ . Since  $L_n^F(\hat{\mathbf{w}}_n) \geq L_n^F(\mathbf{w}_n^F)$  and (A.3)-(A.5) hold under the assumptions of Theorem 2, it suffices for (2.26) to show that

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{\hat{L}_n^*(\mathbf{w}) - L_n^*(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1), \quad (\text{A.8})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{\hat{A}_n(\mathbf{w}) - A_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1), \quad (\text{A.9})$$

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{\hat{B}_n(\mathbf{w}) - B_n(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1), \quad (\text{A.10})$$

and

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \left| \frac{L_n^F(\mathbf{w}) - L_n^*(\mathbf{w})}{R_n^*(\mathbf{w})} \right| = o_p(1). \quad (\text{A.11})$$

To prove (A.9), note first that

$$\begin{aligned} & \left| \hat{A}_n(\mathbf{w}) - A_n(\mathbf{w}) \right| \leq \left| \boldsymbol{\mu}'_n(I - P^*(\mathbf{w}))'(\Sigma_n^{-1} - \hat{\Sigma}_n^{-1})\mathbf{e}_n \right| \\ & + \left| \boldsymbol{\mu}'_n(\hat{P}^*(\mathbf{w}) - P^*(\mathbf{w}))'(\hat{\Sigma}_n^{-1} - \Sigma_n^{-1})\mathbf{e}_n \right| + \left| \boldsymbol{\mu}'_n(\hat{P}^*(\mathbf{w}) - P^*(\mathbf{w}))'\Sigma_n^{-1}\mathbf{e}_n \right| \\ & \equiv (1) + (2) + (3). \end{aligned} \quad (\text{A.12})$$

Assumption 2 implies

$$\sup_{n \geq 1} \|\Sigma_n^{-1}\| < \infty \text{ and } \sup_{n \geq 1} \|\Sigma_n\| < \infty, \quad (\text{A.13})$$

which, together with (2.24), gives

$$\begin{aligned} (1) & \leq C \|\Sigma_n^{-1} - \hat{\Sigma}_n^{-1}\| \|\mathbf{e}_n\| \|\Sigma_n^{-1/2}(I - P^*(\mathbf{w}))\boldsymbol{\mu}_n\| \\ & = O_p(b_n) R_n^{*1/2}(\mathbf{w}), \end{aligned}$$

where the  $O_p(b_n)$  term is independent of  $\mathbf{w}$ . In view of this, (A.7) and (2.25), one obtains

$$\begin{aligned} & \sup_{\mathbf{w} \in \mathcal{H}_N} \frac{\left| \boldsymbol{\mu}'_n(I - P^*(\mathbf{w}))'(\Sigma_n^{-1} - \hat{\Sigma}_n^{-1})\mathbf{e}_n \right|}{R_n^*(\mathbf{w})} \\ & = \max_{1 \leq l \leq N} \max_{1 \leq j_1 < \dots < j_l \leq M} \sup_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} \frac{\left| \boldsymbol{\mu}'_n(I - P^*(\mathbf{w}))'(\Sigma_n^{-1} - \hat{\Sigma}_n^{-1})\mathbf{e}_n \right|}{R_n^*(\mathbf{w})} \\ & = O_p(b_n) \frac{1}{k_n^{*1/2}} = o_p(1). \end{aligned} \quad (\text{A.14})$$

Let  $A_m = X'_m \Sigma_n^{-1} X_m$  and  $\hat{A}_m = X'_m \hat{\Sigma}_n^{-1} X_m$ . Then, straightforward calculations yield for any  $\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}$  with  $1 \leq j_1 < \dots < j_l \leq M$  and  $1 \leq l \leq N$ ,

$$(3) \leq \sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n \left( \hat{\Sigma}_n^{-1} X_m (\hat{A}_m^{-1} - A_m^{-1}) X'_m + (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) X_m A_m^{-1} X'_m \right) \Sigma_n^{-1} \mathbf{e}_n \right|. \quad (\text{A.15})$$

It follows from (A.13) and (2.24) that

$$\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) X_m A_m^{-1} X'_m \Sigma_n^{-1} \mathbf{e}_n \right| = O_p(b_n) \sum_{m \in \{j_1, \dots, j_l\}} \left\| P_m \Sigma_n^{-1/2} \mathbf{e}_n \right\|.$$

In addition, (A.7) and the first moment bound theorem of Findley and Wei (1993) imply that for  $S_1 = S/2 > N/\theta$ ,

$$\begin{aligned} & P_x \left( \max_{1 \leq l \leq N} \max_{1 \leq j_1 < \dots < j_l \leq M} \sup_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left\| P_m \Sigma_n^{-1/2} \mathbf{e}_n \right\|}{R_n^*(\mathbf{w})} > k_n^{*-1/2+\theta} \right) \\ & \leq C \cdot k_n^{*S_1/2} k_n^{*-S_1\theta} \sum_{l=1}^N \sum_{j_l=l}^M \dots \sum_{j_1=1}^{j_2-1} \frac{l}{D_n^{S_1/2}(j_l)} \\ & \leq C \cdot k_n^{*S_1/2} k_n^{*-S_1\theta} k_n^{*-S_1/2+N}. \end{aligned}$$

Combining the above two equations with (2.25) and the dominated convergence theorem, we get

$$\max_{1 \leq l \leq N} \max_{1 \leq j_1 < \dots < j_l \leq M} \sup_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n (\hat{\Sigma}_n^{-1} - \Sigma^{-1}) X_m A^{-1} X'_m \Sigma_n^{-1} \mathbf{e}_n \right|}{R_n^*(\mathbf{w})} = o_p(1). \quad (\text{A.16})$$

Some algebraic manipulations yield

$$\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n \hat{\Sigma}_n^{-1} X_m (\hat{A}_m^{-1} - A_m^{-1}) X'_m \Sigma_n^{-1} \mathbf{e}_n \right| = O_p(b_n) \sum_{m \in \{j_1, \dots, j_l\}} \left\| P_m \Sigma_n^{-1/2} \mathbf{e}_n \right\|.$$

Therefore, by an argument similar to that used to prove (A.16),

$$\max_{1 \leq l \leq N} \max_{1 \leq j_1 < \dots < j_l \leq M} \sup_{\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}} \frac{\sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n \hat{\Sigma}_n^{-1} X_m (\hat{A}_m^{-1} - A_m^{-1}) X'_m \Sigma_n^{-1} \mathbf{e}_n \right|}{R_n^*(\mathbf{w})} = o_p(1). \quad (\text{A.17})$$

We conclude from (A.15), (A.16) and (A.17) that

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \frac{\left| \boldsymbol{\mu}'_n (\hat{P}^*(\mathbf{w}) - P^*(\mathbf{w}))' \Sigma_n^{-1} \mathbf{e}_n \right|}{R_n^*(\mathbf{w})} = o_p(1). \quad (\text{A.18})$$

Finally, straightforward calculations and (2.24) yield that for any  $\mathbf{w} \in \mathcal{H}_{j_1, \dots, j_l}$  with  $1 \leq j_1 < \dots < j_l \leq M$  and  $1 \leq l \leq N$ ,

$$\begin{aligned} (2) &\leq \sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n \hat{\Sigma}_n^{-1} X_m (\hat{A}_m^{-1} - A_m^{-1}) X'_m (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) \mathbf{e}_n \right| \\ &+ \sum_{m \in \{j_1, \dots, j_l\}} \left| \boldsymbol{\mu}'_n (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) X_m (\hat{A}_m^{-1} - A_m^{-1}) X'_m (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) \mathbf{e}_n \right| \\ &= O_p(b_n^2). \end{aligned}$$

This and (2.25) imply

$$\sup_{\mathbf{w} \in \mathcal{H}_N} \frac{\left| \boldsymbol{\mu}'_n (\hat{P}^*(\mathbf{w}) - P^*(\mathbf{w}))' (\hat{\Sigma}_n^{-1} - \Sigma_n^{-1}) \mathbf{e}_n \right|}{R_n^*(\mathbf{w})} = O_p(b_n^2/k_n^*) = o_p(1). \quad (\text{A.19})$$

Now the desired conclusion (A.9) follows from (A.12), (A.14), (A.18) and (A.19). The proofs of (A.8), (A.10) and (A.11) are similar to that of (A.9). The details are thus skipped.

Before proving Theorem 3, we need an auxiliary lemma.

**Lemma 2.** Assume (2.12), (2.14) and (3.6). Then for any  $1 \leq q \leq n-1$ ,

$$\|\Sigma_n^{-1}(q) - \Sigma_n^{-1}\| \leq C \sqrt{\sum_{j \geq q+1} |a_j| \sum_{j \geq q+1} j |a_j|}, \quad (\text{A.20})$$

where  $a_j$ 's are defined as in (3.1).

*Proof.* It follows from (2.12), (2.14), (3.6) and Theorem 3.8.4 of Brillinger (1975) that

$$\sum_{j \geq 1} j|a_j| < \infty. \quad (\text{A.21})$$

In view of (3.2) and (3.3), one has

$$\begin{aligned} \|\Sigma_n^{-1}(q) - \Sigma_n^{-1}\| &\leq \|\mathbf{T}_n - \mathbf{T}_n(q)\| \|\mathbf{D}_n^{-1}\| \|\mathbf{T}_n\| + \|\mathbf{T}_n(q)\| \|\mathbf{D}_n^{-1} - \mathbf{D}_n^{-1}(q)\| \|\mathbf{T}_n\| \\ &+ \|\mathbf{T}_n(q)\| \|\mathbf{D}_n^{-1}(q)\| \|\mathbf{T}_n - \mathbf{T}_n(q)\|. \end{aligned} \quad (\text{A.22})$$

It is easy to see that

$$\|\mathbf{D}_n^{-1}(q)\| \leq C \text{ and } \|\mathbf{D}_n^{-1}\| \leq C. \quad (\text{A.23})$$

Moreover, by (A.13) and (A.23)

$$\|\mathbf{T}_n\| \leq \|\mathbf{T}_n' \mathbf{D}_n^{-1} \mathbf{T}_n\| \|\mathbf{D}_n\| = \gamma_0 \|\Sigma_n^{-1}\| \leq C, \quad (\text{A.24})$$

and

$$\|\mathbf{D}_n^{-1}(q) - \mathbf{D}_n^{-1}\| \leq \|\mathbf{D}_n^{-1}\| \|\mathbf{D}_n^{-1}(q)\| \|\mathbf{D}_n(q) - \mathbf{D}_n\| \leq C(\sigma_q^2 - \sigma_\alpha^2) \leq C \sum_{j \geq q+1} a_j^2. \quad (\text{A.25})$$

According to (A.21)-(A.25), it remains to prove that

$$\|\mathbf{T}_n(q) - \mathbf{T}_n\|^2 \leq C \sum_{j \geq q+1} |a_j| \sum_{j \geq q+1} j|a_j|. \quad (\text{A.26})$$

By making use of Theorem 2.2 of Baxter (1962), it can be shown that

$$\|\mathbf{T}_n(q) - \mathbf{T}_n\|_\infty \leq C \sum_{j \geq q+1} |a_j|, \quad (\text{A.27})$$

and

$$\|\mathbf{T}_n(q) - \mathbf{T}_n\|_1 \leq C \sum_{j \geq q+1} j|a_j|, \quad (\text{A.28})$$

where for an  $m \times n$  matrix  $\mathbf{B} = (b_{ij})_{1 \leq i \leq m, 1 \leq j \leq n}$ ,  $\|\mathbf{B}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |b_{ij}|$  and  $\|\mathbf{B}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |b_{ij}|$ . The desired conclusion (A.26) now follows from (A.27), (A.28) and  $\|\mathbf{T}_n(q) - \mathbf{T}_n\|^2 \leq \|\mathbf{T}_n(q) - \mathbf{T}_n\|_1 \|\mathbf{T}_n(q) - \mathbf{T}_n\|_\infty$ .

*Proof of Theorem 3.* We will first show that for each  $1 \leq k \leq q_n$

$$\mathbb{E} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(k) \hat{\mathbf{e}}_{t+1,k} \right\|^r \leq C \left( \frac{k}{n} \right)^{r/2}, \quad (\text{A.29})$$

where  $N_0 = n - q_n$ ,  $\hat{\mathbf{e}}_t(k) = (\hat{e}_t, \dots, \hat{e}_{t-k+1})'$  and  $\hat{e}_{t+1,k} = \hat{e}_{t+1} + \mathbf{a}'(k) \hat{\mathbf{e}}_t(k)$  with  $\mathbf{a}'(k) = (a_1(k), \dots, a_k(k))$ . Define  $\mathbf{e}_t(k) = (e_t, \dots, e_{t-k+1})'$ ,  $e_{t+1,k} = e_{t+1} + \mathbf{a}'(k) \mathbf{e}_t(k)$  and  $\mathbf{z}_n = (z_1, \dots, z_n)' = (I - H_{d_n}) \mathbf{w}(d_n)$ , where  $\mathbf{w}(d_n) = (w_1(d_n), \dots, w_n(d_n))' = (\sum_{j=d_n}^\infty \theta_j x_{1j}, \dots, \sum_{j=d_n+1}^\infty \theta_j x_{nj})'$ . Let  $\{\mathbf{o}_i = (o_{1i}, \dots, o_{ni})', 1 \leq$



$i \leq d_n\}$  be an orthonormal basis of the column space of  $X(d_n)$ . Then, it holds that  $H_{d_n} = \sum_{i=1}^{d_n} \mathbf{o}_i \mathbf{o}_i'$ . Moreover, one has for  $1 \leq k \leq q_n$ ,

$$\begin{aligned}
& \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(k) \hat{e}_{t+1,k} \\
&= \frac{1}{N_0} \sum_{t=q_n}^{n-1} \left\{ \mathbf{e}_t(k) + \mathbf{z}_t(k) - \sum_{i=1}^{d_n} v_i \mathbf{o}_t^{(i)}(k) \right\} \left\{ e_{t+1,k} + z_{t+1,k} - \sum_{i=1}^{d_n} v_i o_{t+1,k}^{(i)} \right\} \\
&= \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(k) e_{t+1,k} + \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{z}_t(k) e_{t+1,k} - \sum_{i=1}^{d_n} v_i \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) e_{t+1,k} \right\} \\
&+ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(k) z_{t+1,k} + \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{z}_t(k) z_{t+1,k} - \sum_{i=1}^{d_n} v_i \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) z_{t+1,k} \right\} \\
&- \sum_{i=1}^{d_n} v_i \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(k) o_{t+1,k}^{(i)} \right\} - \sum_{i=1}^{d_n} v_i \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{z}_t(k) o_{t+1,k}^{(i)} \right\} \\
&+ \sum_{i=1}^{d_n} \sum_{j=1}^{d_n} v_i v_j \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) o_{t+1,k}^{(j)} \right\} \\
&:= (\text{I}) + \dots + (\text{IX}), \tag{A.30}
\end{aligned}$$

where  $\mathbf{z}_t(k) = (z_t, \dots, z_{t-k+1})'$ ,  $z_{t+1,k} = z_{t+1} + \mathbf{a}'(k) \mathbf{z}_t(k)$ ,  $v_i = \mathbf{o}_t' \mathbf{e}$ ,  $\mathbf{o}_t^{(i)}(k) = (o_{ti}, \dots, o_{t-k+1,i})'$  and  $o_{t+1,k}^{(i)} = o_{t+1,i} + \mathbf{a}'(k) \mathbf{o}_t^{(i)}(k)$ . By Lemmas 3 and 4 of Ing and Wei (2003),

$$\mathbb{E} \|(\text{I})\|^r \leq C \left( \frac{k}{n} \right)^{r/2}. \tag{A.31}$$

Theorem 2.2 of Baxter (1962) and (3.6) ensure that the spectral density of  $e_{t+1,k}$  is bounded above, and hence by (3.5), Lemma 2 of Wei (1987) and Minkowski's Inequality,

$$\begin{aligned}
\mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} z_{t-l} e_{t+1,k} \right|^r &\leq C n^{-r} \mathbb{E} \left( \sum_{t=q_n}^{n-1} z_{t-l}^2 \right)^{r/2} \\
&\leq C n^{-r/2} \mathbb{E} \left( \frac{1}{n} \sum_{t=1}^n w_t^2(d_n) \right)^{r/2} \leq C n^{-r/2} \left( \sum_{j>d_n} |\theta_j| \right)^r,
\end{aligned}$$

for all  $0 \leq l \leq k-1$ . As a result,

$$\mathbb{E} \|(\text{II})\|^r \leq C \left( \frac{k}{n} \right)^{r/2} \left( \sum_{j>d_n} |\theta_j| \right)^r. \tag{A.32}$$

By (A.13), (3.5), the boundedness of the spectral density of  $e_{t+1,k}$ , and Lemma 2 of Wei (1987), one has for all  $1 \leq i \leq d_n$  and all  $0 \leq l \leq k-1$ ,

$$\mathbb{E} |v_i|^{2r} \leq C \mathbb{E} \left( \sum_{t=1}^n o_{ti}^2 \right)^r \leq C, \tag{A.33}$$

and

$$\mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} o_{t-l,i} e_{t+1,k} \right|^{2r} \leq C n^{-2r} \mathbb{E} \left( \sum_{t=1}^n o_{ti}^2 \right)^r \leq C n^{-2r}. \quad (\text{A.34})$$

Making use of (A.33), (A.34), the convexity of  $x^r, x \geq 0$ , and the Cauchy-Schwarz inequality, we obtain

$$\begin{aligned} \mathbb{E} \|(\text{III})\|^r &\leq \mathbb{E} \left( \sum_{i=1}^{d_n} |v_i| \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) e_{t+1,k} \right\| \right)^r \\ &\leq d_n^r d_n^{-1} \sum_{i=1}^{d_n} \mathbb{E} \left( |v_i|^r \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) e_{t+1,k} \right\|^r \right) \\ &\leq C d_n^r d_n^{-1} \sum_{i=1}^{d_n} \left( \mathbb{E} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) e_{t+1,k} \right\|^{2r} \right)^{1/2} \\ &\leq C \frac{d_n^r k^{r/2}}{n^r}. \end{aligned} \quad (\text{A.35})$$

Following an argument similar to that used to prove (A.32), we have for  $0 \leq l \leq k-1$ ,

$$\begin{aligned} \mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} z_{t+1,k} e_{t-l} \right|^r &\leq C n^{-r} \mathbb{E} \left( \sum_{t=q_n}^{n-1} z_{t+1,k}^2 \right)^{r/2} \\ &\leq C n^{-r} \mathbb{E} \left\{ \left( 1 + \sum_{j=1}^k |a_j(k)|^2 \right) \left( \sum_{t=1}^n w_t^2(d_n) \right) \right\}^{r/2} \\ &\leq C n^{-r/2} \left( \sum_{j>d_n} |\theta_j| \right)^r, \end{aligned}$$

and hence

$$\mathbb{E} \|(\text{IV})\|^r \leq C \left( \frac{k}{n} \right)^{r/2} \left( \sum_{j>d_n} |\theta_j| \right)^r. \quad (\text{A.36})$$

Similarly, for all  $0 \leq l \leq k-1$ ,

$$\mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} z_{t-l} z_{t+1,k} \right|^r \leq C \mathbb{E} \left( \frac{1}{N_0} \sum_{t=1}^n w_t^2(d_n) \right)^r \leq C \left( \sum_{j>d_n} |\theta_j| \right)^{2r},$$

yielding

$$\mathbb{E} \|(\text{V})\|^r \leq C \left( \frac{k}{n} \right)^{r/2} \left( n^{1/4} \sum_{j>d_n} |\theta_j| \right)^{2r}. \quad (\text{A.37})$$

By an argument analogous to (A.35), it holds that

$$\mathbb{E}\|(\text{VII})\|^r \leq C \frac{d_n^r k^{r/2}}{n^r}. \quad (\text{A.38})$$

According to (A.33),

$$\begin{aligned} \mathbb{E}\|(\text{VI})\|^r &\leq C d_n^r d_n^{-1} \sum_{i=1}^{d_n} \left\{ \mathbb{E}|v_i|^r \mathbb{E} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) z_{t+1,k} \right\|^r \right\} \\ &\leq C d_n^r d_n^{-1} \sum_{i=1}^{d_n} \mathbb{E} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{o}_t^{(i)}(k) z_{t+1,k} \right\|^r. \end{aligned}$$

Moreover, Minkowski's inequality and the Cauchy-Schwarz inequality yield that for all  $1 \leq i \leq d_n$  and  $0 \leq l \leq k-1$ ,

$$\begin{aligned} \mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} o_{t-l,i} z_{t+1,k} \right|^r &\leq \mathbb{E} \left\{ \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} o_{t-l,i}^2 \right)^{r/2} \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} Z_{t+1,k}^2 \right)^{r/2} \right\} \\ &\leq C n^{-r/2} \mathbb{E} \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} z_{t+1,k}^2 \right)^{r/2} \leq C n^{-r/2} \left( \sum_{j>d_n} |\theta_j| \right)^r \end{aligned}$$

As a result,

$$\mathbb{E}\|(\text{VI})\|^r \leq C \frac{k^{r/2}}{n^{r/2}} \left( d_n \sum_{j>d_n} |\theta_j| \right)^r. \quad (\text{A.39})$$

Similarly, it can be shown that

$$\mathbb{E}\|(\text{VIII})\|^r \leq C \frac{k^{r/2}}{n^{r/2}} \left( d_n \sum_{j>d_n} |\theta_j| \right)^r, \quad (\text{A.40})$$

and

$$\mathbb{E}\|(\text{IX})\|^r \leq C \frac{d_n^{2r} k^{r/2}}{n^r}. \quad (\text{A.41})$$

Consequently, (A.29) follows from (A.30)-(A.32) and (A.35)-(A.41).

Let

$$G_n = \max_{1 \leq k \leq q_n} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(k) \hat{e}_{t+1,k} \right\|.$$

Then, for any  $M > 0$ , one obtains from (A.29) and Chebyshev's inequality that

$$\begin{aligned} P \left( G_n > M \frac{q_n^{1/2+1/r}}{n^{1/2}} \right) &= P \left( G_n^r > M^r \left( \frac{q_n}{n} \right)^{r/2} q_n \right) \\ &\leq \frac{C}{q_n^{r/2+1} M^r} \sum_{k=1}^{q_n} k^{r/2} \leq \frac{C}{M^r}. \end{aligned}$$

Hence

$$\max_{1 \leq k \leq q_n} \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(k) \hat{e}_{t+1,k} \right\| = O_p \left( \frac{q_n^{1/2+1/r}}{n^{1/2}} \right). \quad (\text{A.42})$$

In the following, we shall show that

$$\left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) - \Sigma_{q_n} \right\| = o_p(1). \quad (\text{A.43})$$

Note first that

$$\begin{aligned} & \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) - \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(q_n) \mathbf{e}'_t(q_n) \right\| \\ & \leq C \left\{ \frac{1}{N_0} \sum_{t=q_n}^{n-1} \|\hat{\mathbf{e}}_t(q_n) - \mathbf{e}_t(q_n)\|^2 \right. \\ & \quad \left. + \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} \|\hat{\mathbf{e}}_t(q_n) - \mathbf{e}_t(q_n)\|^2 \right)^{1/2} \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} \|\mathbf{e}_t(q_n)\|^2 \right)^{1/2} \right\}. \end{aligned}$$

Straightforward calculations imply

$$\frac{1}{N_0} \sum_{t=q_n}^{n-1} \|\hat{\mathbf{e}}_t(q_n) - \mathbf{e}_t(q_n)\|^2 \leq C \frac{q_n}{N_0} \|\hat{\mathbf{e}}_n - \mathbf{e}_n\|^2 \leq \frac{C q_n}{N_0} \left\{ \mathbf{e}'_n H_{d_n} \mathbf{e}_n + \|\mathbf{w}_n(d_n)\|^2 \right\},$$

$E(\mathbf{e}'_n H_{d_n} \mathbf{e}_n) \leq C d_n$ ,  $E(\|\mathbf{w}_n(d_n)\|^2) \leq n C (\sum_{j>d_n} |\theta_j|)^2$ , and  $E(N_0^{-1} \sum_{t=q_n}^{n-1} \|\mathbf{e}_t(q_n)\|^2) \leq C q_n$ . As a result,

$$\begin{aligned} & \left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) - \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(q_n) \mathbf{e}'_t(q_n) \right\| \\ & = O_p \left( \frac{q_n d_n}{n} + q_n \left( \sum_{j>d_n} |\theta_j| \right)^2 + \frac{q_n d_n^{1/2}}{n^{1/2}} + q_n \sum_{j>d_n} |\theta_j| \right) = o_p(1). \end{aligned} \quad (\text{A.44})$$

Moreover, Lemma 2 of Ing and Wei (2003) yields

$$\left\| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \mathbf{e}_t(q_n) \mathbf{e}'_t(q_n) - \Sigma_{q_n} \right\| = O_p \left( \frac{q_n}{n^{1/2}} \right) = o_p(1). \quad (\text{A.45})$$

Combining (A.44) and (A.45) leads to the desired conclusion (A.43).

By making use of (A.42) and (A.43), we next show that

$$\left\| \hat{\mathbf{T}}_n(q_n) - \mathbf{T}_n(q_n) \right\| = O_p \left( \frac{q_n^{1+1/r}}{n^{1/2}} \right). \quad (\text{A.46})$$

It follows from (A.43) and (A.13) that

$$\lim_{n \rightarrow \infty} P(Q_n) \equiv \lim_{n \rightarrow \infty} P((N_0^{-1} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n))^{-1} \text{ exists}) = 1, \quad (\text{A.47})$$

and

$$\left\| \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) \right)^{-1} \right\| I_{Q_n} = O_p(1). \quad (\text{A.48})$$

In addition, an argument given in Proposition 3.1 of Ing, Chiou and Guo (2013) implies

$$\left\| \hat{\mathbf{T}}_n(q_n) - \mathbf{T}_n(q_n) \right\| \leq C q_n^{1/2} \max_{1 \leq k \leq q_n} \|\hat{\mathbf{a}}(k) - \mathbf{a}(k)\|, \quad (\text{A.49})$$

where  $\hat{\mathbf{a}}(k) = (\hat{a}_1(k), \dots, a_k(k))'$ . Since on  $Q_n$ ,

$$\max_{1 \leq k \leq q_n} \|\hat{\mathbf{a}}(k) - \mathbf{a}(k)\| \leq \left\| \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) \right)^{-1} \right\| G_n, \quad (\text{A.50})$$

(A.46) is ensured by (A.42) and (A.47)-(A.50).

The proof of (3.11) is also reliant on

$$\|\hat{\mathbf{D}}_n^{-1}(q_n) - \mathbf{D}_n^{-1}(q_n)\| = O_p \left( \frac{q_n^{1/r}}{n^{1/2}} + \frac{q_n^{1+(2/r)}}{n} \right), \quad (\text{A.51})$$

which is in turn implied by (A.23) and

$$\|\hat{\mathbf{D}}_n(q_n) - \mathbf{D}_n(q_n)\| = O_p \left( \frac{q_n^{1/r}}{n^{1/2}} + \frac{q_n^{1+(2/r)}}{n} \right). \quad (\text{A.52})$$

To prove (A.52), note first that on the set  $Q_n$ ,

$$\begin{aligned} & \max_{1 \leq k \leq q_n} |\hat{\sigma}_k^2 - \sigma_k^2| \\ & \leq \max_{1 \leq k \leq q_n} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{e}_{t+1,k}^2 - \sigma_k^2 \right| + \left\| \left( \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{\mathbf{e}}_t(q_n) \hat{\mathbf{e}}'_t(q_n) \right)^{-1} \right\| G_n^2. \end{aligned} \quad (\text{A.53})$$

Moreover, by (3.8), (3.9), Lemma 6 of Ing and Wei (2005) and an argument similar to that used to prove (A.29), it holds that for all  $1 \leq k \leq q_n$ ,

$$\mathbb{E} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{e}_{t+1,k}^2 - \sigma_k^2 \right|^r \leq C n^{-r/2},$$

and hence

$$\max_{1 \leq k \leq q_n} \left| \frac{1}{N_0} \sum_{t=q_n}^{n-1} \hat{e}_{t+1,k}^2 - \sigma_k^2 \right| = O_p \left( \frac{q_n^{1/r}}{n^{1/2}} \right). \quad (\text{A.54})$$

Similarly, we have

$$|\hat{\gamma}_0 - \gamma_0| = O_p(n^{-1/2}). \quad (\text{A.55})$$

Combining (3.10), (A.42), (A.47), (A.48), (A.53)-(A.55) and

$$\|\hat{\mathbf{D}}_n(q_n) - \mathbf{D}_n(q_n)\| = \max\{|\hat{\gamma}_0 - \gamma_0|, \max_{1 \leq k \leq q_n} |\hat{\sigma}_k^2 - \sigma_k^2|\}$$

yields the desired conclusion (A.52). The proof is completed by noticing that (3.11) is an immediate consequence of (A.20), (A.46) and (A.51).

## ACKNOWLEDGMENTS

The research of Ching-Kang Ing was supported in part by the Academia Sinica Investigator Award, and that of Shu-Hui Yu was partially supported by the National Science Council of Taiwan under grant NSC 99-2118-M-390-002. We would like to thank the editors and two anonymous referees for their insightful and constructive comments, which greatly improve the presentation of this paper.

## REFERENCES

- H. Akaike (1974). *A new look at the statistical model identification*. *IEEE Trans. Automatic Control* **19** 716-723.
- T. Ando and K.-C. Li (2014). A model-averaging approach for high-dimensional regression. *J. Amer. Statist. Assoc.* forthcoming.
- D. W. K. Andrews (1991). Asymptotic optimality of generalized  $C_L$ , cross-validation, and generalized cross-validation in regression with heteroskedastic errors. *J. Economet.* **4** 359-377.
- G. Baxter (1962). An Asymptotic Result for the Finite Predictor. *Math. Scand.* **10** 137-144.
- D. R. Brillinger (1975). *Time Series: Data Analysis and Theory*. Holt, Rinehart and Winston, New York.
- N. H. Chan, S.-F. Huang and C.-K. Ing (2013). Moment bound and mean squared prediction errors of long-memory time series. *Ann. Statist.* **41** 1268-1298.
- D. F. Findley and C. Z. Wei (1993). Moment bounds for deriving time series CLT's and model selection procedures. *Statist. Sinica* **3** 453-470.
- B. E. Hansen (2007). Least squares model averaging. *Econometrica* **75** 1175-1189.
- B. E. Hansen and J. S. Racine (2012). Jackknife model averaging. *J. Economet.* **167** 38-46.
- C.-K. Ing (2007). Accumulated prediction errors, information criteria and optimal forecasting for autoregressive time series. *Ann. Statist.* **35** 1238-1277.
- C.-K. Ing, H. T. Chiou and M. H. Guo (2013). Estimation of inverse autocovariance matrices for long memory processes. Technical Report.

- C.-K. Ing and T. L. Lai (2011). A stepwise regression method and consistent model selection for high-dimensional sparse linear models. *Statist. Sinica* **21** 1473–1513.
- C.-K. Ing and C.-Z. Wei (2003). On same-realization prediction in an infinite-order autoregressive process. *J. Multivariate Anal.* **85** 130–155.
- C.-K. Ing and C.-Z. Wei (2005). Order selection for same-realization predictions in autoregressive processes. *Ann. Statist.* **33** 2423–2474.
- T. L. Lai and C.-Z. Wei (1982). Least squares estimates in stochastic regression models with applications to identification and control of dynamic systems. *Ann. Statist.* **1** 154–166.
- S. Lee and A. Karagrigoriou (2001). An asymptotically optimal selection of the order of a linear process. *Sankhya A* **63** 93–106.
- G. Leung and A. R. Barron (2006). Information theory and mixing least-squares regressions. *IEEE Trans. Information Theory* **52** 3396–3410.
- K.-C. Li (1987). Asymptotic optimality for  $C_p$ ,  $C_L$ , cross-validation and generalized cross-validation: discrete index set. *Ann. Statist.* **15** 958–975.
- Q. Liu and R. Okui (2013). Heteroscedasticity-robust  $C_p$  Model Averaging. *Economet. J.* **16** 463–472.
- Q. Liu, R. Okui, and A. Yoshimura (2013). Generalized Least Squares Model Averaging. Technical Report.
- C. L. Mallows (1973). Some comments on  $C_p$ . *Technometrics* **15** 661–675.
- T. L. McMurphy and D. N. Politis (2010). Banded and tapered estimates for autocovariance matrices and the linear process bootstrap. *J. Time Series Anal.* **31** 471–482.
- J. Rissanen (1986). Order estimation by accumulated prediction errors. In *Essays in Time Series and Allied Processes* (J. Gani and M. P. Priestley, eds.) *J. Appl. Probab.* **23A** 55–61.
- R. Shibata (1981). An optimal selection of regression variables. *Biometrika* **68** 45–54.
- V. N. Temlyakov (2000). Weak greedy algorithms. *Adv. Comput. Math.* **12** 213–227.
- A. T. K. Wan, X. Zhang, and G. Zou (2010). Least squares model averaging by Mallows criterion. *J. Economet.* **156** 227–283.
- Z. Wang, S. Paterlini, F. Gao, and Y. Yang (2014). Adaptive minimax regression estimation over sparse hulls. *J. Machine Learning Research* **15** 1675–1711.
- C. Z. Wei (1987). Adaptive prediction by least squares predictors in stochastic regression models with applications to time series. *Ann. Statist.* **15** 1667–1682.
- X. Wei and Y. Yang (2012). Robust forecast combinations. *J. Economet.* **166** 224–236.
- P. Whittle (1960). Bounds for the moments of linear and quadratic forms in independent variables. *Theory Probab. Appl.* **5** 302–305.



- W. B. Wu and M. Pourahmadi (2003). Nonparametric estimation of large covariance matrices of longitudinal data. *Biometrika* **90** 831–844.
- W. B. Wu and M. Pourahmadi (2009). Banding sample covariance matrices of stationary processes. *Statist. Sinica* **19** 1755–1768.
- Y. Yang (2001). Adaptive regression by mixing. *J. Amer. Statist. Assoc.* **96** 574–586.
- Y. Yang (2007). Prediction/Estimation with Simple Linear Model: Is It Really that Simple? *Economet. Theory* **23** 1–36.
- S. H. Yu, C.-C. Lin and H.-W. Cheng (2012). A note on mean squared prediction error under the unit root model with deterministic trend. *J. Time Series Anal.* **33** 276–286.
- Z. Yuan and Y. Yang (2005). Combining linear regression model: when and how? *J. Amer. Statist. Assoc.* **100** 1202–1214.
- X. Zhang, A. T. K. Wan and G. Zou (2013). Model averaging by jackknife criterion in models with dependent data. *J. Economet.* **174** 82–94.
- A. Zygmund (1959). Trigonometric Series, 2nd ed. Cambridge Univ. Press.